

# UC Irvine

## UC Irvine Previously Published Works

### Title

A Bayesian phylogenetic hidden Markov model for B cell receptor sequence analysis.

### Permalink

<https://escholarship.org/uc/item/62j344xw>

### Journal

PLoS computational biology, 16(8)

### ISSN

1553-734X

### Authors

Dhar, Amrit  
Ralph, Duncan K  
Minin, Vladimir N  
[et al.](#)

### Publication Date

2020-08-01

### DOI

10.1371/journal.pcbi.1008030

### Copyright Information

This work is made available under the terms of a Creative Commons Attribution License, available at <https://creativecommons.org/licenses/by/4.0/>

Peer reviewed

# A Bayesian Phylogenetic Hidden Markov Model for B Cell Receptor Sequence Analysis

Amrit Dhar<sup>1,2</sup>, Duncan K. Ralph<sup>2</sup>, Vladimir N. Minin<sup>3,\*</sup>, Frederick A. Matsen IV<sup>2,\*</sup>

<sup>1</sup>Department of Statistics, University of Washington, Seattle

<sup>2</sup>Fred Hutchinson Cancer Research Center

<sup>3</sup>Department of Statistics, University of California, Irvine

\*corresponding authors: [vminin@uci.edu](mailto:vminin@uci.edu), [matsen@fredhutch.org](mailto:matsen@fredhutch.org)

July 1, 2019

## Abstract

The human body is able to generate a diverse set of high affinity antibodies, the soluble form of B cell receptors (BCRs), that bind to and neutralize invading pathogens. The natural development of BCRs must be understood in order to design vaccines for highly mutable pathogens such as influenza and HIV. BCR diversity is induced by naturally occurring combinatorial “V(D)J” rearrangement, mutation, and selection processes. Most current methods for BCR sequence analysis focus on separately modeling the above processes. Statistical phylogenetic methods are often used to model the mutational dynamics of BCR sequence data, but these techniques do not consider all the complexities associated with B cell diversification such as the V(D)J rearrangement process. In particular, standard phylogenetic approaches assume the DNA bases of the progenitor (or “naive”) sequence arise independently and according to the same distribution, ignoring the complexities of V(D)J rearrangement. In this paper, we introduce a novel approach to Bayesian phylogenetic inference for BCR sequences that is based on a phylogenetic hidden Markov model (phylo-HMM). This technique not only integrates a naive rearrangement model with a phylogenetic model for BCR sequence evolution but also naturally accounts for uncertainty in all unobserved variables, including the phylogenetic tree, via posterior distribution sampling.

## Introduction

One of the most important features of the adaptive immune system is its ability to create a wide variety of high affinity antibodies, the soluble form of B cell receptors (BCRs), that bind to and neutralize pathogens in the body. The initial BCR diversity is generated by randomly joining together various gene segments in a process called V(D)J rearrangement; after an initial testing process the cells reach the “naive” state. When stimulated by binding to foreign material called “antigen,” B cells diversify further by entering germinal centers (GCs) in the secondary lymphoid organs and going through an affinity maturation process. During the GC reaction, B cells mutate rapidly in a process called somatic hypermutation (SHM), and the high affinity clones are positively selected for via clonal expansion. We would like to better understand the GC mutation and selection processes, because insight into mutational pathways from naive to mature BCR sequences could aid in the development of vaccines for highly mutable pathogens such as influenza and HIV ([Mascola and Haynes, 2013](#)). We have developed a new statistical inference framework that better estimates these mutational pathways and quantifies uncertainty in these estimates.

Rational vaccine design efforts depend on accurate inference of full evolutionary paths from a given naive sequence to the corresponding mature BCR sequences in a GC. By understanding the mutational pathways that lead to broadly neutralizing antibodies (bNAbs), vaccines could then be constructed that induce the production of these bNAbs in the body ([Stamatatos et al., 2017](#)). For instance, in the case of HIV, most bNAbs are generally not observed until after a long period of chronic infection, requiring prospective studies to characterize them ([Liao et al., 2013](#)). While there have been many such studies that experimentally test longitudinal B cell samples spanning from an initial HIV infection to the development of mature bNAbs ([Liao et al., 2013](#); [Doria-Rose et al., 2014, 2016](#)), obtaining early prospective

samples is difficult. One way to avoid the process described above is to infer intermediate lineage sequences computationally, synthesize them in the laboratory, and test their binding and neutralization abilities (Doria-Rose et al., 2014; Simonich et al., 2019). Consequently, we will focus on inferring mutational pathways through ancestral sequence reconstruction, a commonly used technique in computational phylogenetics.

Much of the existing BCR sequence analysis literature focuses on modeling either the V(D)J recombination process, or the phylogenetic diversification process, but not both. Elhanati et al. (2015) develop a likelihood-based model that encompasses the V(D)J rearrangement, SHM, and clonal selection processes; however, they implicitly assume clonal sequences arise independently and do not consider the phylogenetic structure of SHM. Hoehn et al. (2017, 2019) introduce novel codon substitution models for ancestral lineage inference that encodes SHM context-dependent mutational effects but does not account for V(D)J rearrangement dynamics in the naive sequence within their maximum likelihood phylogenetic inference framework. Yaari et al. (2013) provide estimates of context-dependent SHM substitution probabilities and motif mutability scores based on an aggregated dataset consisting of the synonymous codon positions of productive antibody sequences. Ralph and Matsen IV (2016a) are able to infer naive BCR sequences using an HMM-based approach that models the V(D)J rearrangement process but assumes independent evolution across the different lineages. While these efforts have contributed greatly to our understanding of GC dynamics, we believe that the performance of clonal lineage and ancestral sequence inference procedures can be enhanced by using an evolutionary model for SHM that also accounts for uncertainty in the naive rearrangement process.

Kepler (2013) has developed a likelihood-based SHM modeling framework that jointly estimates the naive sequence and the associated clonal tree and incorporates information about the V(D)J rearrangement process. However, this work does not consider phylogenetic or ancestral sequence uncertainty in the naive sequence estimation procedure. While there is evidence to suggest that ancestral sequence estimation is robust to phylogenetic uncertainty in other settings (Hanson-Smith et al., 2010), the parameter regime of BCR diversification is quite different from that of this previous work, leaving open the question of whether incorporating phylogenetic uncertainty would aid in ancestral sequence inference. Therefore, we would like to construct a phylogenetic inference procedure that not only allows for easy quantification of phylogenetic and ancestral sequence uncertainty but also models the V(D)J recombination as an informative prior for the naive sequence at the root of a phylogenetic tree describing the evolution of one clonal lineage.

In this paper, we propose a Bayesian approach to phylogenetic inference for clonal sequences that is based on a phylogenetic hidden Markov model (phylo-HMM) (Siepel and Haussler, 2005). Our phylo-HMM models both the naive rearrangement and SHM processes. The Bayesian framework allows us to naturally account for uncertainty in all unobserved variables, including a phylogenetic tree, via posterior distribution sampling. We perform simulation-based experiments to show that naive sequence and phylogenetic inference performed jointly provides higher-quality estimates than those obtained by considering these inferences separately. Our application to real data reveals significant uncertainty in naive and ancestral sequences, confirming the importance of a Bayesian approach.

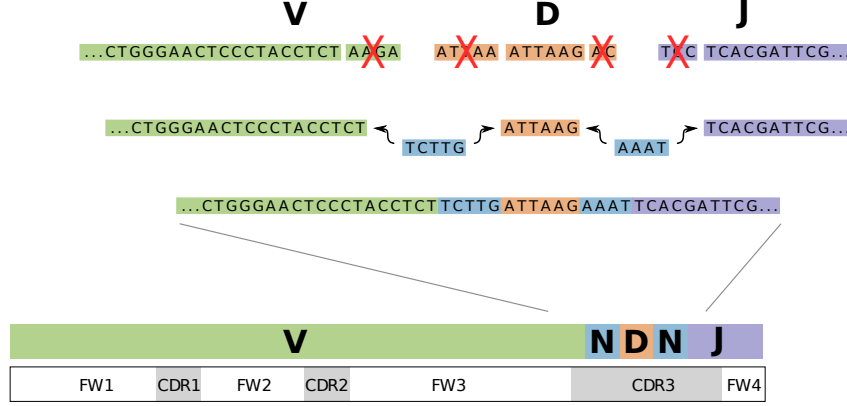
## Methods

### Overview

The immune system is able to generate a diverse set of naive BCR sequences due to the V(D)J rearrangement process (Figure 1a). In this process, the B cell first randomly selects V, D, and J gene segments (i.e. DNA sequences) from the respective gene pools in the body. Before joining the gene segments together, the B cell randomly deletes nucleotides at both ends of the V-D and D-J junction regions and randomly inserts nucleotides in the same junction regions (i.e. non-templated insertions). BCRs comprise both heavy chain and light chain sequences, where the former goes through the VDJ rearrangement process described above while the latter undergoes VJ recombination, a process similar to VDJ recombination without D germline genes; in this paper, we focus solely on heavy chain inference, but our methods and implementation extend to light chain inference too. Although the V(D)J rearrangement process samples germline genes from the same gene pools for all the naive BCRs in a given individual, different people may have different collections of germline genes (Watson et al., 2017). BCR sequences can be partitioned into framework (FWK) and complementarity-determining (CDR) regions. The BCR binding affinity is largely determined by the sequence segments in the CDR regions, and among all the CDR regions the CDR3 region contributes the most to antigen-binding specificity and has the highest amount of sequence

variability. The **partis** software program (Ralph and Matsen IV, 2016a,b) treats the naive sequence generative process as a discrete-time Markov chain (DTMC) going from left to right across the sequence bases and also permits maximum likelihood inference of the associated model parameters, which are defined according to the V(D)J rearrangement dynamics. We emphasize that the “time” of this DTMC represents position along the sequence, in contrast to the continuous time Markov chain described below modeling the substitution process through chronological time.

(A)



(B)

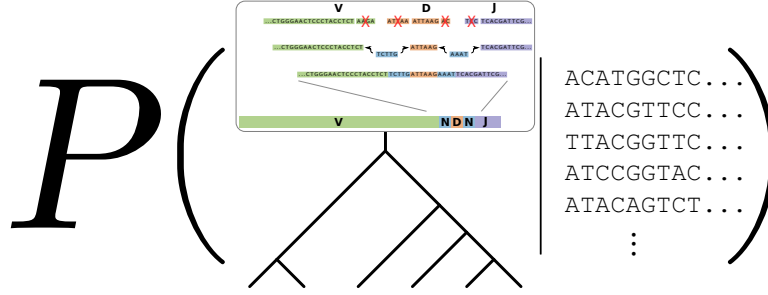


Figure 1: (A) A schematic representation of the naive rearrangement process. First, V (green), D (orange), and J (purple) genes are randomly selected from the respective gene pools in the body. Then, nucleotides are randomly deleted (red X’s) from both ends of the V-D and D-J junction regions and random bases (blue) are added to the same junction regions before the V, D, and J germline genes can be joined together. The BCR sequences can be partitioned into framework (FWK) and complementarity-determining (CDR) regions. This image was taken from (Ralph and Matsen IV, 2016a). (B) Our Bayesian phylo-HMM jointly models V(D)J recombination at the root of the tree (using an HMM) and then subsequent diversification (via a phylogenetic tree). We do posterior inference conditioning on the observed sequence alignment in a clonal family, but not on a fixed inferred naive sequence.

As we attempt to characterize the affinity maturation process in GCs, we use the concept of a clonal family (CF) to help analyze BCR repertoire datasets that result from high-throughput sequencing experiments. A clonal family represents a set of BCR sequences that originate from the same naive rearrangement event; in practice, it is simpler to define these families as groups of BCR sequences that share the same naive sequence (Ralph and Matsen IV, 2016b). In this work, we use the latter definition of a CF. The CF definition used here relies on the basic assumption that the unmutated common ancestor of a collection of clonally related sequences can only be identified by the corresponding DNA sequence. There is a chance that two different naive B cells with identical BCR sequences could seed multiple GCs and form their own lineages, but the observed clones from the two GCs would be collapsed into a single CF because both lineages share the same naive sequence. Thus, a correctly inferred CF under this definition will be a cluster of sequences that derive from naive B cells with the same BCR sequence; because GC mutation and selection occur at the sequence level, all sequences in a CF go through the same mutation and selection processes. Using this CF definition, Ralph and Matsen IV (2016b) describe how to cluster BCR sequences from large-scale repertoire datasets into CFs with high accuracy using **partis**. We want to emphasize that repertoire datasets can be clustered and pre-processed in many different

ways, but **partis** provides a convenient way to cluster repertoire sequences and produce CF-specific multiple sequence alignments (MSAs).

As mentioned in the opening section, naive sequences accumulate mutations via SHM and the corresponding B cells undergo cellular replication. Phylogenetic tree models provide a realistic and mathematically convenient way to represent the aforementioned B cell evolutionary dynamics of a CF. In particular, these models define a likelihood at each MSA position/site as a function of unknown parameters. These parameters consist of a tree topology, branch lengths, and continuous-time Markov chain (CTMC) substitution model parameters. While B cell evolution in a CF is not independent across the different site positions, site-specific phylogenetic likelihoods provide a convenient first approximation in modeling this phenomenon. To help us illustrate how these phylogenetic models work, we provide an example tree visualization for 4 sequences (Figure S1).

Given a tree topology with branch lengths, we use a CTMC substitution model to calculate the probabilities of state changes along the branches of the tree. Specifically if  $t$  denotes a branch length on a tree, CTMC substitution models allow one to calculate  $p_{ij}(t)$ , which denotes the probability of going from state  $i$  to state  $j$  on a branch of length  $t$ , where  $i, j \in \{A, G, C, T\}$ . It is common to use a reversible CTMC substitution model on a tree (Felsenstein, 2004); a reversible substitution model is a Markov substitution model that, if started at stationarity, can be run backwards in time, with the resulting backward Markov model following the same probability law as the original forward model. The standard phylogenetic generative process can be described at each alignment site as follows: 1) a DNA state at the root node is drawn independently according to the same 4-state discrete distribution and 2) the states at the other nodes are sampled in a pre-order traversal using the computed CTMC probabilities at each branch  $p_{ij}(t)$ ; this model is a special case of a directed graphical model (Lauritzen, 1996) that probabilistically generates sequence alignments.

Clearly, standard phylogenetic models do not account for naive rearrangement dynamics at the root because, as we discussed above, the root state at each sequence position is sampled independently according to an identical distribution. Instead, if we draw root sequence states from the DTMC mentioned earlier, we would obtain a sequence evolution model that more accurately describes B cell evolutionary dynamics. Thus, we formulate our phylo-HMM to consist of a hidden state DTMC model for naive sequences that explicitly incorporates V(D)J rearrangement information, and an emission distribution that generates sequence alignments conditional on the naive sequence that is based on phylogenetic likelihoods. This phylo-HMM hopefully leads to more accurate naive sequence estimates and, as a result of that, higher-quality intermediate ancestral sequence estimates. We introduce a pictorial representation of our Bayesian phylo-HMM to make clear our target of inference (Figure 1b).

## Notation and Assumptions

We now introduce some notation and assumptions that will be used throughout this paper. Let  $\mathbf{D} = \{D_i^{(j)}\}_{i=1:m, j=1:n}$  denote the MSA of  $m$  clonal DNA sequences of length  $n$ . We define  $\mathbf{Y}_{\text{naive}} = \{Y_{\text{naive}}^{(j)}\}_{j=1:n}$  and  $\hat{\mathbf{Y}}_{\text{naive}} = \{\hat{Y}_{\text{naive}}^{(j)}\}_{j=1:n}$  to be the corresponding length- $n$  naive sequence random variable and point estimate, respectively. We let  $\tau$  represent a tree topology with  $m$  tips and a root branch length; in total, this topology has  $m$  internal nodes and  $2m - 1$  branch lengths. We assume that the ancestral sequence at the root of  $\tau$  is  $\mathbf{Y}_{\text{naive}}$ . Furthermore, we define  $\mathbf{t} = \{t_i\}_{i=1:(2m-1)}$  to be the branch lengths associated with  $\tau$ . Let  $\mathbf{Y}_{\text{int}} = \{Y_i^{(j)}\}_{i=1:(m-1), j=1:n}$  denote the internal nodes of  $\tau$  excluding the naive sequence  $\mathbf{Y}_{\text{naive}}$ . For convenience, we let  $\mathbf{D}^{(j)} = \{D_i^{(j)}\}_{i=1:m}$  and  $\mathbf{Y}_{\text{int}}^{(j)} = \{Y_i^{(j)}\}_{i=1:(m-1)}$  symbolize the observed sequence data and unobserved ancestral sequence data, respectively, at MSA site  $j \in \{1, \dots, n\}$ . Conditioned on the root sequence  $\mathbf{Y}_{\text{naive}}$ , we assume that the ancestral states at each site in the MSA evolve independently along the phylogeny  $\tau$  according to a general time-reversible (GTR) substitution model (Tavaré, 1986). Let  $\mathbf{e} = \{e_{AC}, e_{AG}, e_{AT}, e_{CG}, e_{CT}, e_{GT}\}$  and  $\boldsymbol{\pi} = \{\pi_A, \pi_C, \pi_G, \pi_T\}$  represent the GTR exchangeability rates and equilibrium base frequencies, respectively. We also account for phylogenetic rate variation among sites by employing a discrete gamma distribution with a fixed number of rate classes  $K$  (Yang, 1994, 1996) and define  $\alpha$  to be the associated gamma shape parameter, denote  $\mathbf{r} = \{r_1(\alpha), \dots, r_K(\alpha)\}$  as the set of discrete rates deterministically induced by  $\alpha$ , and let  $\mathbf{r}^* = \{r_{(j)}^*\}_{j=1:n}$  represent the discrete rates chosen at each site in the MSA. In theory, we would like to compute phylogenetic likelihoods with branch lengths scaled by  $r_{(j)}^*$  and mix over  $r_{(j)}^* \sim \text{Gamma}(\alpha, \alpha)$  for  $j \in \{1, \dots, n\}$ . However, these integrals are generally intractable so Yang (1994) suggested dividing the  $\text{Gamma}(\alpha, \alpha)$  distribution into  $K$  equal-probability rate classes, with the mean rate in each class used to represent all rates in that class. In practice, we use the  $\text{Categorical}(\mathbf{r}, \mathbf{p})$  distribution to define

the models on  $r_{(j)}^*$  for  $j \in \{1, \dots, n\}$ , where  $\mathbf{p}$  is a, possibly unnormalized, probability vector; in addition, this model is used to represent more general discrete distributions.

From a statistical point of view, it is common to assume the naive sequence root node is a leaf node holding naive sequence bases connected to a “virtual root node” (i.e. what we call “root” node for our phylogenetic model described in the previous subsection) via a branch length of 0. Even though it may seem like we have described a rooted tree model above, it turns out that under the assumptions of a reversible substitution model and a nucleotide distribution at the virtual root starting at stationarity, the Pulley Principle, first discussed in (Felsenstein, 1981), states that the virtual root may be placed anywhere on the tree without affecting the likelihood. This implies that the model described above does not correspond to a single rooted tree, but an equivalence class of rooted trees that maps to a unique unrooted tree. This is an important distinction as our phylo-HMM will in fact use an unrooted tree model, which will be justified when we describe our posterior sampling approach.

## Phylo-HMM Description

Phylo-HMMs are special cases of directed graphical models (Lauritzen, 1996) and treat evolution as a combination of two Markov processes: one across the sites in the MSA and one down the phylogeny. They are commonly used for sequence-level segmentation problems such as gene prediction and detection of highly-conserved regions (Siepel and Haussler, 2005). In fact, a phylo-HMM is similar in structure to a standard HMM; the main difference between the two model classes is that a phylo-HMM uses a phylogenetic likelihood as its emission probability distribution, while standard HMMs usually specify simpler emission distributions. Our BCR-specific phylo-HMM specifies a Markov process along  $\mathbf{Y}_{\text{naive}}$  and, conditional on  $\mathbf{Y}_{\text{naive}}$ , a phylogenetic evolutionary process down the given tree. Phylo-HMMs have not been applied to BCR sequence analysis before and we believe this biologically realistic probabilistic model is uniquely suited to provide higher-quality naive sequence and ancestral sequence estimates compared to those obtained under current state-of-the-art methods.

To help us describe the phylo-HMM generative process, we provide an illustration of the associated graphical model diagram for an example alignment with  $m = 3$  sequences and  $n = 3$  sites (Figure 2). The naive sequence “hidden state” prior distribution  $p(\mathbf{Y}_{\text{naive}})$  decomposes to  $p(Y_{\text{naive}}^{(1)}) \prod_{j=2}^n p(Y_{\text{naive}}^{(j)} | Y_{\text{naive}}^{(j-1)})$  and the bases are generated sequentially; these prior probabilities depend on hyperparameters that can be set using the `partis` software package (Ralph and Matsen IV, 2016a,b). For the tree topology  $\tau$ , we assume that a tree is drawn from the Uniform distribution over  $(m + 1)$ -tip unrooted trees; this seems like a strange choice given that we described  $\tau$  as a rooted topology above, but this decision will be justified when we discuss how to perform Bayesian inference under the phylo-HMM. The branch lengths  $\mathbf{t}$  and the gamma shape parameter  $\alpha$  are assumed to be *a priori* independent and to follow Exponential( $\lambda$ ) distributions, where  $\lambda$  is some prespecified rate. The GTR exchangeability rates  $\mathbf{e}$  and equilibrium base frequencies  $\boldsymbol{\pi}$  are usually assumed to come from six-dimensional and four-dimensional Dirichlet distributions, respectively.

For each MSA site  $j \in \{1, \dots, n\}$ ,  $r_{(j)}^*$  *a priori* follows the Categorical( $\mathbf{r}, (\frac{1}{K}, \dots, \frac{1}{K})$ ) distribution. Then, at each site  $j \in \{1, \dots, n\}$ , we assume that  $\mathbf{Y}_{\text{int}}^{(j)}$  and  $\mathbf{D}^{(j)}$  are generated by drawing DNA states from CTMC transition probability matrices based on augmented branch lengths  $\mathbf{t} \times r_{(j)}^*$ . For example, in Figure 2, we first sample  $Y_1^{(j)}$  from  $p(Y_1^{(j)} | Y_{\text{naive}}^{(j)})$ , which is a row vector distribution in the CTMC transition probability matrix for the “branch length”  $t_1 \times r_{(j)}^*$ , where  $j \in \{1, \dots, n\}$ . Once we have sampled  $Y_1^{(j)}$  for  $j = 1, \dots, n$ , we can draw  $Y_2^{(j)}$  and  $D_3^{(j)}$  using similar row vector distributions from CTMC transition probability matrices for the “branch lengths”  $t_2 \times r_{(j)}^*$  and  $t_3 \times r_{(j)}^*$ , respectively. We can recursively continue this process until we generate states at  $D_1^{(j)}$  and  $D_2^{(j)}$  for  $j \in \{1, \dots, n\}$ .

## Posterior Distribution Inference

We are interested in sampling from the posterior distribution  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} | \mathbf{D})$ :

$$\begin{aligned} & p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} | \mathbf{D}) \\ & \propto p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}}, \mathbf{D}) \\ & = p(\mathbf{r}^*, \mathbf{Y}_{\text{int}}, \mathbf{D} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}) p(\tau) p(\mathbf{t}) p(\boldsymbol{\pi}) p(\mathbf{e}) p(\alpha) p(\mathbf{r} | \alpha) p(\mathbf{Y}_{\text{naive}}) \end{aligned}$$

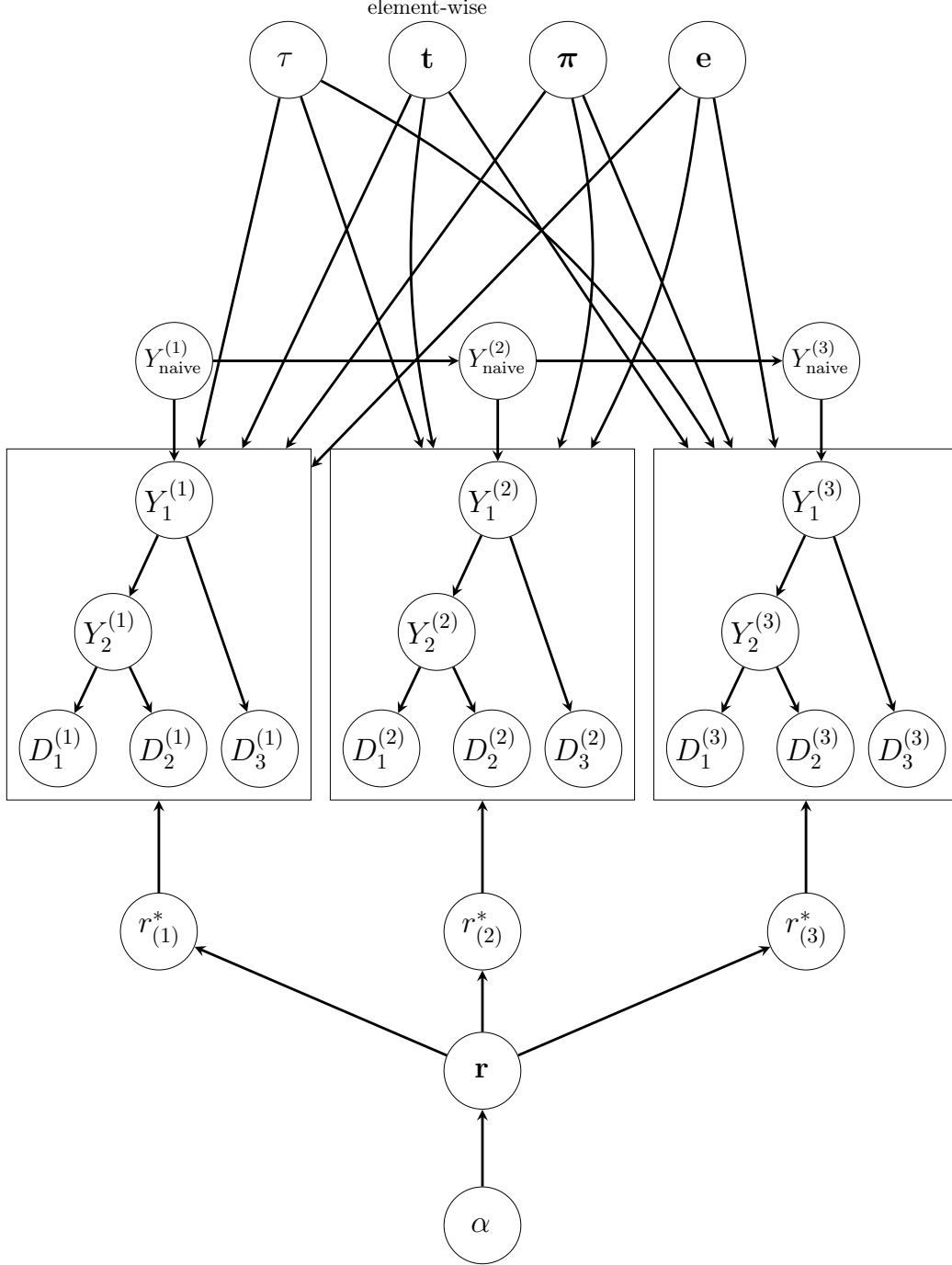


Figure 2: The phylo-HMM graphical model diagram for an example alignment with  $m = 3$  sequences and  $n = 3$  sites. The  $\tau$ ,  $\mathbf{t}$ ,  $\boldsymbol{\pi}$ , and  $\mathbf{e}$  nodes represent the 4-tip unrooted tree topology, the associated 5 branch lengths, the GTR exchangeability rates, and GTR equilibrium base frequencies, respectively. The parameter  $\alpha$  denotes the gamma shape parameter associated with the  $K$ -class discrete gamma distribution, which is used to model phylogenetic rate variation among sites;  $\mathbf{r}$  symbolizes the vector of  $K$  discrete rates that is deterministically induced by  $\alpha$ . The set of nodes  $\mathbf{r}^* = \{r_{(1)}^*, r_{(2)}^*, r_{(3)}^*\}$  defines the rates that are drawn from  $\mathbf{r}$  at each particular site. The  $\mathbf{Y}_{\text{naive}} = \{Y_{\text{naive}}^{(1)}, Y_{\text{naive}}^{(2)}, Y_{\text{naive}}^{(3)}\}$  “hidden state” node collection represents the Markov process that stochastically generates the naive sequence in our phylo-HMM. The node sets  $\{Y_i^{(j)}\}_{i=1:2, j=1:3}$  and  $\mathbf{D} = \{D_i^{(j)}\}_{i=1:3, j=1:3}$  denote the internal nodes of  $\tau$  excluding the naive sequence  $\mathbf{Y}_{\text{naive}}$  and the observed MSA, respectively. We draw plates around the  $\mathbf{Y}_{\text{int}}^{(j)}$  and  $\mathbf{D}^{(j)}$  node sets for  $j \in \{1, 2, 3\}$  to indicate that any directed edges touching a plate apply to all nodes in the plate (except for edges that originate from  $\mathbf{t}$ , which apply element-wise to the nodes in the plate).



$$= \left\{ \prod_{j=1}^n p(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, r_{(j)}^*, \mathbf{Y}_{\text{int}}^{(j)}) p(\mathbf{Y}_{\text{int}}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, r_{(j)}^*, \mathbf{Y}_{\text{naive}}^{(j)}) p(r_{(j)}^* \mid \mathbf{r}) \right\} \\ \times p(\tau) p(\mathbf{t}) p(\boldsymbol{\pi}) p(\mathbf{e}) p(\alpha) p(\mathbf{r} \mid \alpha) p(\mathbf{Y}_{\text{naive}}),$$

where this model decomposition results from the definition of a directed graphical model. We also factorize the posterior distribution in the following way:

$$p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} \mid \mathbf{D}) \\ = p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}) \\ \times p(\mathbf{Y}_{\text{naive}} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{D}) \\ \times p(\mathbf{r}^*, \mathbf{Y}_{\text{int}} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}, \mathbf{D}).$$

This formulation is useful because it suggests that we can generate draws from the posterior distribution by sampling sequentially from three conditional probability distributions. Conceptually, to sample from the posterior, we have to draw in-order: 1) the phylogeny-related parameters, 2) the naive sequence, and 3) the ancestral sequences. We describe how to perform these three sampling steps in the following subsections and provide a complete summary of the sampling process in Algorithm 1.

### Tree Sampling

Our strategy for sampling from  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$  is to first draw a large pool of observations from an easy-to-sample proposal distribution  $q$  and then perform weighted bootstrap with weights  $\frac{p}{q}$  on those samples to obtain approximate draws from the correct distribution. This “sampling-importance-resampling” (SIR) algorithm is a sample filtering method that finds use in a wide variety of statistical applications (Smith and Gelfand, 1992; Gordon et al., 1993; Andrieu et al., 2010). The original SIR algorithm resampled observations with replacement, but there are theoretical and practical considerations that make resampling without replacement more attractive (Skare et al., 2003; Gelman et al., 2013). A thorough review of SIR sampling can be found in (Gelman and Meng, 2004, Chapter 24). We use the SIR algorithm to sample from  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$  because we want to take advantage of already-existing software programs for Bayesian phylogenetic inference while incorporating biologically-realistic details into our phylo-HMM.

Our phylogeny proposal distribution  $q$  comes from the **RevBayes** software package (Höhna et al., 2016) because this is a proposal that is close to the target distribution  $p$  while also being easy to sample from. In short, this  $q$  is traditional Bayesian phylogenetic analysis with a point estimate of the naive sequence. In more detail, we first input the MSA  $\mathbf{D}$  into the **partis** package, obtain a naive sequence point estimate  $\hat{\mathbf{Y}}_{\text{naive}}^{\text{partis}}$ , and create an augmented MSA  $\mathbf{D}^* = \{\mathbf{D}, \hat{\mathbf{Y}}_{\text{naive}}^{\text{partis}}\}$ . This allows us to generate Markov chain Monte Carlo (MCMC) samples from  $q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$  and provides a convenient way to sample trees that have  $m$  tips and an informative root branch length emanating from the naive sequence. In our **RevBayes** analysis, we require that the prior components of  $q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$  be defined as we discussed in the previous subsection for  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$  and as we will see, this assumption is critical to the validity of our technique.

For the purposes of this discussion, let us assume that we sample a large number (say  $N_{\text{pool}}$ ) of parameter values from  $q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$  in **RevBayes**. As we briefly mentioned above, we resample  $N_{\text{final}}$  times without replacement from the set of  $N_{\text{pool}}$   $q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$  draws. Each  $q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$  sample is assigned a sampling weight  $w = p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}) / q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$ . While it seems odd to use the ratio of posterior probabilities that are conditional on different datasets as bootstrap weights, the only technical requirement on  $p$  and  $q$  is that the parameters of interest (i.e.  $\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha$ , and  $\mathbf{r}$ ) have the same support, which is indeed the case in our situation. Smith and Gelfand (1992) suggest picking  $N_{\text{pool}}$  and  $N_{\text{final}}$  so  $\frac{N_{\text{pool}}}{N_{\text{final}}} \geq 10$  while Rubin (1987) proposed a safe rule-of-thumb to be  $\frac{N_{\text{pool}}}{N_{\text{final}}} = 20$ ; we use  $\frac{N_{\text{pool}}}{N_{\text{final}}} = 20$  in all the applied experiments conducted in this paper. It may not be immediately clear how one would efficiently compute the numerator in  $w$  so we express  $w$  in the following form:

$$w = \frac{p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})}{q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)} \\ = \frac{p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{D}) / p(\mathbf{D})}{q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{D}^*) / q(\mathbf{D}^*)} \\ = \frac{p(\mathbf{D} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})}{q(\mathbf{D}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})} \frac{q(\mathbf{D}^*)}{p(\mathbf{D})}$$



---

**Algorithm 1:** Posterior Sampling of  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} \mid \mathbf{D})$ 


---

**Input:** CF multiple sequence alignment  $\mathbf{D}$ , number of discrete rates  $K$

$N_{\text{pool}}, N_{\text{final}}$  ( $N_{\text{pool}}/N_{\text{final}} \approx 20$ )

**Output:**  $N_{\text{final}}$  samples of  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} \mid \mathbf{D})$

*Tree Sampling* —  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$

1. Run **partis** on input data  $\mathbf{D}$ .

$\Rightarrow \mathbf{D}^* = \{\mathbf{D}, \hat{\mathbf{Y}}_{\text{naive}}^{\text{partis}}\}, \hat{p}(\mathbf{Y}_{\text{naive}})$

2. Run **RevBayes** MCMC on the augmented MSA  $\mathbf{D}^*$ .

$\Rightarrow N_{\text{pool}}$  samples of  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}) \sim q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)$

3. Run the SIR algorithm without replacement on the  $N_{\text{pool}}$   $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})$  proposal samples with weights  $w = \frac{p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})}{q(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D}^*)}$ .

$\Rightarrow N_{\text{final}}$  samples of  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$

*Naive Sequence Sampling* —  $p(\mathbf{Y}_{\text{naive}} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{D})$

For each sample  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} \mid \mathbf{D})$ :

For each site  $j \in \{n, \dots, 1\}$ :

1. Draw  $Y_{\text{naive}}^{(j)}$  using our phylo-HMM-based backward sampling procedure.

$\Rightarrow N_{\text{final}}$  samples of  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}} \mid \mathbf{D})$

*Intermediate Ancestral Sequence Sampling* —  $p(\mathbf{r}^*, \mathbf{Y}_{\text{int}} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}, \mathbf{D})$

For each sample  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}} \mid \mathbf{D})$ :

For each site  $j \in \{1, \dots, n\}$ :

1. Sample  $r_{(j)}^*$  according to probabilities proportional to  $p(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*)$ .

2. Sample  $\mathbf{Y}_{\text{int}}^{(j)}$  in a pre-order fashion using the standard ASR distribution at internal nodes on a  $r_{(j)}^*$ -scaled tree.

$\Rightarrow N_{\text{final}}$  samples of  $(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}}) \sim p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} \mid \mathbf{D})$

---

$$\begin{aligned}
&= \frac{\sum_{\mathbf{Y}_{\text{naive}}} p(\mathbf{D} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}) p(\mathbf{Y}_{\text{naive}} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}) q(\mathbf{D}^*)}{q(\mathbf{D}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})} \frac{q(\mathbf{D}^*)}{p(\mathbf{D})} \\
&= \frac{\sum_{\mathbf{Y}_{\text{naive}}} p(\mathbf{D} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, \mathbf{Y}_{\text{naive}}) p(\mathbf{Y}_{\text{naive}}) q(\mathbf{D}^*)}{q(\mathbf{D}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})} \frac{q(\mathbf{D}^*)}{p(\mathbf{D})} \\
&= \frac{\sum_{\mathbf{Y}_{\text{naive}}} \left\{ p(Y_{\text{naive}}^{(1)}) \prod_{j=2}^n p(Y_{\text{naive}}^{(j)} \mid Y_{\text{naive}}^{(j-1)}) \prod_{k=1}^n p(\mathbf{D}^{(k)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, Y_{\text{naive}}^{(k)}) \right\} q(\mathbf{D}^*)}{q(\mathbf{D}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})} \frac{q(\mathbf{D}^*)}{p(\mathbf{D})},
\end{aligned}$$

where the transition from the second line to the third line is due to the priors on  $\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha$ , and  $\mathbf{r}$  for both  $p$  and  $q$  being equal and the decomposition between the fourth and sixth lines is a result of  $d$ -separation conditional independencies (Lauritzen, 1996) as follows. Specifically,  $\mathbf{Y}_{\text{naive}} \perp \{\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}\}$  because every undirected path in the graphical model between  $\mathbf{Y}_{\text{naive}}$  and  $\{\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}\}$  does not contain any “conditioned” child nodes and  $\mathbf{D} \perp \alpha \mid \{\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, \mathbf{Y}_{\text{naive}}\}$  as all paths between  $\mathbf{D}$  and  $\alpha$  are “blocked” by the intermediate node  $\mathbf{r}$ . The final denominator term above  $q(\mathbf{D}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r})$  is the usual phylogenetic likelihood, here calculated by **RevBayes**. Note that the marginal likelihood ratio

$\frac{q(\mathbf{D}^*)}{p(\mathbf{D})}$  does not affect the bootstrap sampling because the sampling probabilities are proportional to  $w$  and the likelihood ratio is a constant with respect to  $\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha$ , and  $\mathbf{r}$  in the  $N_{\text{pool}}$  sampling weights. The numerator term in the final equation for  $w$  looks complicated, but is actually the phylo-HMM likelihood with the “hidden state” probabilities  $p(Y_{\text{naive}}^{(1)})$  and  $p(Y_{\text{naive}}^{(j)} | Y_{\text{naive}}^{(j-1)})$  for  $j \in \{2, \dots, n\}$  and the “emission probabilities”  $p(\mathbf{D}^{(k)} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, Y_{\text{naive}}^{(k)})$  that marginalize over the site-wise rates  $r_{(k)}^*$  for  $k \in \{1, \dots, n\}$ . We can efficiently calculate the phylo-HMM likelihood using the forward algorithm (Rabiner, 1986), but this approach requires us to be able to compute the phylo-HMM emission probabilities in a straightforward manner. The computation of the emission probabilities is of interest because the hidden state probabilities in  $p(\mathbf{Y}_{\text{naive}})$  can be easily inferred using **partis**, which we now call  $\hat{p}(\mathbf{Y}_{\text{naive}})$ .

It turns out that we can, again, leverage existing software tools to help us efficiently compute the emission probabilities  $p(\mathbf{D}^{(k)} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, Y_{\text{naive}}^{(k)})$  for  $k \in \{1, \dots, n\}$ . The key point is to recognize that  $p(\mathbf{D}^{(k)} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, Y_{\text{naive}}^{(k)})$  is an entry in the Felsenstein likelihood vector at the  $Y_{\text{naive}}^{(k)}$  node, which denotes the probability of the observed data at only the tips that descend from node  $Y_{\text{naive}}^{(k)}$ , given the conditioned state of node  $Y_{\text{naive}}^{(k)}$ . These vectors are commonly used within a post-order tree traversal algorithm to compute standard phylogenetic likelihoods (Felsenstein, 1981). Let  $\mathbf{F}_u = (F_{u1}, \dots, F_{um})^T$  be the vector of partial likelihoods at node  $u$ , where  $F_{ui}$  denotes the probability of the observed data at only the tips that descend from node  $u$ , given that the state of node  $u$  is  $i$ . Because we utilize an unrooted tree model in our phylo-HMM, we can use a Pulley Principle argument (Felsenstein, 1981) to show that a standard phylogenetic likelihood on our tree can be represented as  $p(\mathbf{D}^{(k)} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \mathbf{r}, Y_{\text{naive}}^{(k)}) \pi_{Y_{\text{naive}}^{(k)}}^{(k)}$ . Thus, we can compute the phylo-HMM emission probabilities by first calculating the standard site-wise phylogenetic likelihoods on the same tree and then dividing out the  $Y_{\text{naive}}^{(k)}$  stationary probabilities. To compute these standard site-specific phylogenetic likelihoods, we make use of the **libpll** C library (Flouri et al., 2014), which is a versatile high-performance software library for phylogenetic analysis. Of course, we could have instead used a rooted tree model in our phylo-HMM, performed our own post-order tree traversal algorithm for likelihood computations, and extracted out the appropriate entries in the site-wise Felsenstein likelihood vectors, but we want our inference technique to scale to the large datasets that result from high-throughput BCR sequencing and no currently-existing software package meets this requirement. In the next subsection, we describe how to generate naive sequence draws given the approximate phylogenetic tree samples from  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} | \mathbf{D})$ .

### Naive Sequence Sampling

To sample naive sequences, we exploit the fact that our phylo-HMM is essentially a standard HMM with a naive-conditional phylogenetic likelihood as its emission distribution. We draw the “hidden state” naive sequence  $\mathbf{Y}_{\text{naive}}$  from  $p(\mathbf{Y}_{\text{naive}} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{D})$  by adapting the hidden state posterior sampling technique for standard HMMs to be used for our specialized phylo-HMM (Scott, 2002). Just as we perform the forward algorithm by recursively computing and caching intermediate phylo-HMM likelihoods (i.e. forward probabilities) going left to right across the MSA, we can sample  $\mathbf{Y}_{\text{naive}}$  by doing a backward pass through the phylo-HMM and drawing the  $Y_{\text{naive}}^{(j)}$  states starting at site  $n$  and ending at the first alignment site. In fact, the maximum a posteriori probability (MAP) estimate of the hidden state sequence  $\mathbf{Y}_{\text{naive}}$  is obtained using a similar procedure called the Viterbi algorithm (Rabiner, 1986; Scott, 2002). This backward sampling procedure can actually use the previously cached forward probabilities in the calculation of the sampling probabilities at each site, which is convenient because we already had to run the forward algorithm to sample the phylogeny-related parameters from  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r} | \mathbf{D})$ . Once  $\mathbf{Y}_{\text{naive}}$  has been sampled, we can proceed to draw the intermediate ancestral states  $\mathbf{Y}_{\text{int}}$  from the conditional distribution  $p(\mathbf{r}^*, \mathbf{Y}_{\text{int}} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}, \mathbf{D})$ .

### Intermediate Ancestral Sequence Sampling

Just as we did for our naive sequence sampling, we sample the intermediate ancestral states  $\mathbf{Y}_{\text{int}}$  by utilizing a modified version of the standard ancestral sequence reconstruction (ASR) technique used in phylogenetics (Nielsen, 2002). It is important to note that sampling from  $p(\mathbf{r}^*, \mathbf{Y}_{\text{int}} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{Y}_{\text{naive}}, \mathbf{D})$  can be reduced to drawing  $(r_{(j)}^*, \mathbf{Y}_{\text{int}}^{(j)})$  pairs from  $p(r_{(j)}^*, \mathbf{Y}_{\text{int}}^{(j)} | \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, \mathbf{D}^{(j)})$  for each MSA site  $j \in \{1, \dots, n\}$ , which is justified by  $d$ -separation conditional independencies (Lauritzen, 1996). At each

site  $j \in \{1, \dots, n\}$ , we first sample  $r_{(j)}^*$  and then  $\mathbf{Y}_{\text{int}}^{(j)}$  according to the previously-described distribution:

$$\begin{aligned} & p\left(r_{(j)}^*, \mathbf{Y}_{\text{int}}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, \mathbf{D}^{(j)}\right) \\ &= p\left(r_{(j)}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, \mathbf{D}^{(j)}\right) p\left(\mathbf{Y}_{\text{int}}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, \mathbf{D}^{(j)}\right), \end{aligned}$$

where the above decomposition is based on the definition of conditional probability. We draw the site-specific rates before the site-wise intermediate ancestral states because conditioning on the rates allows for simpler and more efficient ancestral state sampling using an already-existing software package.

It turns out that we can draw  $r_{(j)}^*$  values from  $\mathbf{r}$  with probabilities proportional to  $p(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*)$ , which makes intuitive sense because it implies rates should be sampled according to the site-specific likelihoods with branch lengths scaled by the corresponding rates. To understand why the above statement holds true, we express  $p(r_{(j)}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, \mathbf{D}^{(j)})$  as follows:

$$\begin{aligned} & p\left(r_{(j)}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, \mathbf{D}^{(j)}\right) \\ & \propto p\left(r_{(j)}^*, \mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}\right) \\ &= p\left(r_{(j)}^* \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}\right) p\left(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right) \\ &= p\left(r_{(j)}^* \mid \mathbf{r}\right) p\left(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right) \\ & \propto p\left(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right), \end{aligned}$$

where the transition from the second line to the third line stems from the definition of conditional probability, the transition from the third to fourth line is a result of  $d$ -separation (Lauritzen, 1996), and the fourth-to-fifth line transition is due to the fact that  $r_{(j)}^* \mid \mathbf{r} \sim \text{Categorical}(\mathbf{r}, (\frac{1}{K}, \dots, \frac{1}{K}))$  by assumption. These site-specific likelihoods are in fact almost identical to the naive-conditional phylogenetic likelihoods discussed in the previous two subsections with the only difference being that we now condition on the site-wise rates instead of marginalizing over them.

To illustrate why sampling  $\mathbf{Y}_{\text{int}}^{(j)}$  from  $p(\mathbf{Y}_{\text{int}}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, \mathbf{D}^{(j)})$  is similar to drawing ASRs according to the procedure outlined by Nielsen (2002), we derive the sampling distribution of  $Y_1^{(j)}$ , the most recent common ancestor of  $\mathbf{D}^{(j)}$ . The distribution  $p(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, \mathbf{D}^{(j)})$  can be decomposed in the following manner:

$$\begin{aligned} & p\left(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, \mathbf{D}^{(j)}\right) \\ & \propto p\left(Y_1^{(j)}, \mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right) \\ &= p\left(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right) p\left(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, Y_1^{(j)}\right) \\ &= p\left(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*\right) p\left(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, r_{(j)}^*, Y_1^{(j)}\right), \end{aligned}$$

where, as before, the transition between the second and third lines are due to the definition of conditional probability and the transition from the third line to the fourth line is a result of  $d$ -separation (Lauritzen, 1996). The first term in the resulting expression above,  $p(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, Y_{\text{naive}}^{(j)}, r_{(j)}^*)$ , is the CTMC transition probability between  $Y_{\text{naive}}^{(j)}$  and  $Y_1^{(j)}$  on a scaled branch length  $t_1 \times r_{(j)}^*$  and the second term,  $p(\mathbf{D}^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, r_{(j)}^*, Y_1^{(j)})$ , is an entry in the Felsenstein likelihood vector at the  $Y_1^{(j)}$  node on the given tree with  $r_{(j)}^*$ -scaled branch lengths. This expression for the sampling distribution of  $Y_1^{(j)}$  corresponds with the standard ASR distribution at internal nodes described in (Nielsen, 2002). Thus, by conceptually treating  $Y_{\text{naive}}^{(j)}$  as a sampled root node and scaling the branch lengths by  $r_{(j)}^*$ , we can draw  $Y_1^{(j)} \sim p(Y_1^{(j)} \mid \tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, Y_{\text{naive}}^{(j)}, r_{(j)}^*, \mathbf{D}^{(j)})$  using equation (10) presented by Nielsen (2002). By recursively using  $d$ -separation arguments (Lauritzen, 1996), we can use similar logic to that described above to sample states at the other internal nodes  $\mathbf{Y}_{\text{int}}^{(j)} \setminus Y_1^{(j)}$  in a pre-order fashion as well. Once we have drawn  $\mathbf{Y}_{\text{int}}^{(j)}$  at each alignment site  $j \in \{1, \dots, n\}$ , our three-stage sampling process of  $p(\tau, \mathbf{t}, \boldsymbol{\pi}, \mathbf{e}, \alpha, \mathbf{r}, \mathbf{r}^*, \mathbf{Y}_{\text{naive}}, \mathbf{Y}_{\text{int}} \mid \mathbf{D})$  is completed.

## Implementation

The entire Bayesian phylo-HMM sampling process is implemented in a pipeline called **linearham**, which is available at <https://github.com/matsengrp/linearham>. We developed our phylo-HMM data structure in C++ and, as mentioned before, used the tree structures from the **libpll** C library (Flouri et al., 2014) to ensure our phylogenetic likelihood computations were as fast as possible. We have provided a Docker container (<https://hub.docker.com/r/matsengrp/linearham>) so users can try out the software without installation, as well as specifying the installation dependencies and provide an example command in a Dockerfile.

Our **linearham** program also provides an interface to **partis** to obtain  $\hat{p}(\mathbf{Y}_{\text{naive}})$ . It accepts a repertoire FASTA file as input and can infer the hidden state transition probabilities assuming either the repertoire needs to be partitioned into individual CFs before phylogenetic inference or the “repertoire” is a single CF already, which is useful for researchers that want to run **linearham** inference on a hand-curated CF. It is also possible to specify an external set of  $p(\mathbf{Y}_{\text{naive}})$  parameters.

In addition, **linearham** summarizes its phylogenetic inference output in a user-friendly format. It provides an output FASTA file that contains each sampled amino acid naive sequence and its associated posterior probability, creates a FASTA-like file that maps each sampled amino acid naive sequence to its corresponding set of DNA naive sequences and posterior probabilities, and generates an amino acid naive sequence posterior probability logo using **WebLogo** (Crooks et al., 2004) to visualize the per-site uncertainties. Furthermore, we provide similarly-styled output files for particular naive-to-tip tree lineages of interest by tabulating the posterior probabilities of sampled naive sequences and intermediate ancestral sequences on the lineage. For naive-to-tip lineage analysis, we also create a posterior probability lineage graphic using **Graphviz** (Gansner and North, 2000) that summarizes the different inferred naive-to-tip sequence trajectories and their relative confidences (Gong et al., 2013).

## Simulation Experiments

In our simulation experiments, we focus on validating the accuracy of the naive sequence and ancestral sequence estimates produced by **linearham**. For naive sequence validations, we compare and contrast the performance between **linearham**, **partis**, and the **ARPP** program (Kepler, 2013). The **partis** package provides maximum likelihood naive sequence estimates (Ralph and Matsen IV, 2016a), but its model assumes a star-tree configuration, which is unrealistic for many CFs going through long periods of SHM and affinity maturation (Liao et al., 2013; Doria-Rose et al., 2014, 2016). The **ARPP** program is also a likelihood-based framework that jointly infers a CF tree and the associated naive sequence using information about the V(D)J rearrangement process, but does not quantify phylogenetic or ancestral sequence uncertainty; we use this program in our validations because it is one of the only other programs that estimates CF phylogenies and naive sequences at the same time. We did attempt to use the newer version of **ARPP** (called **Cloanalyst**), but ran into difficulties with the program crashing and were unable to successfully resolve this issue when we contacted the program author.

Similarly, for our ancestral sequence validations, we restrict our comparisons to the **linearham**, **RevBayes**, and **dnaml** (Felsenstein, 2005) packages. As discussed before, the **RevBayes** program performs Bayesian phylogenetic inference on a given MSA, but in this case, we sample CF trees and ASRs from **RevBayes** using an augmented CF sequence alignment that contains the **partis**-inferred naive sequence (i.e.  $\mathbf{D}^*$  from above). This approach to ASR was used in (Simonich et al., 2019) and is similar to that of **linearham** with the main difference being that **RevBayes** ASR sampling is conditional on the **partis**-inferred naive sequence whereas **linearham** ASR inference conditions on naive sequences drawn from a posterior distribution. For all the experiments conducted in this section, we run **RevBayes** using 50,000 MCMC iterations, sampling every 10 iterations, discard the first 500 **RevBayes** samples as burn-in, and sample without replacement 225 times from the 4,500 effective **RevBayes** samples in the case of **linearham** inference. The **dnaml** package performs maximum likelihood phylogenetic inference and generates an ASR conditional on this likelihood-based tree estimate. While it is possible to sample ASRs on a maximum likelihood tree (Nielsen, 2002), **dnaml** only reports the most probable ancestral sequences. We run **dnaml** on an augmented CF sequence alignment that contains the **ARPP** maximum likelihood estimate of the naive sequence. In the following subsections, we describe our simulation experiments in more detail from the data-generating mechanism to the validation results.

## Simulation Setup

To simulate tree topologies with a fixed number of tips in our experiments, we used the single-parameter beta-splitting generative process (Aldous, 1996). The beta-splitting process is able to generate a wide variety of tree topologies ranging from balanced topologies (i.e. trees with approximately equal root-to-tip distances) to imbalanced topologies (i.e. trees with highly variable root-to-tip distances) by varying the associated “balance” parameter  $\beta$ . Intuitively, this process can be seen as a recursive partitioning procedure that, at each tree split, partitions the  $\text{Uniform}(0, 1)$  random numbers corresponding to the “tips” on the current sub-interval according to a  $\text{Beta}(\beta + 1, \beta + 1)$  distribution. As  $\beta \rightarrow \infty$ , the generated trees get closer and closer to balanced binary trees and, as  $\beta \rightarrow -2$ , the simulated topologies look more and more “comb-like” (i.e. imbalanced) (Aldous, 1996).

We are motivated to use this topology-generating process because the level of balance of the tree determines the extent to which a phylogenetic approach to naive sequence estimation improves over a star-tree model. Informally speaking, a phylogenetic approach weights the information coming from tips close to the root (in the imbalanced case) more strongly than tips more distant from the root. Thus we expect a phylogenetic approach to be superior in the imbalanced case. On the other hand, **partis** assumes evolution occurs according to a star tree, which implies the expected number of substitutions from the root to each of the tip sequences should be approximately equal. Thus, for imbalanced trees, that have a large variance in the root-to-tip branch length distances, we would expect **partis** to provide poor naive sequence estimates for sequence datasets generated on those trees compared to a phylogenetic approach. Throughout the rest of this section, we define “tree imbalance” to be the standard deviation of a tree’s root-to-tip distances.

To generate branch lengths for our simulated trees that preserve the shapes induced by the beta-splitting topology prior, we independently draw values from a  $\text{Uniform}(0, 2M)$  distribution, where  $M$  is a constant derived from HIV bnAb lineage trees. Specifically, we ran the PC64 (Landais et al., 2017) and VRC01 (Wu et al., 2015) datasets through **partis** to obtain augmented CF sequence alignments with the **partis**-inferred naive sequence, inferred approximate maximum likelihood CF trees using **FastTree** (Price et al., 2009, 2010), and set  $M$  to be the average across all estimated branch lengths in the two trees ( $\approx 0.0179$ ). We describe these two datasets in more detail in the next section when we discuss our real-world dataset applications. In the subsequent parts of this section, we refer to the inferred **FastTree** trees described above as the PC64 and VRC01 phylogenies, respectively. We emphasize that these trees are only used as the basis for a simulation study and their level of accuracy is not especially important.

Each simulated phylogeny also receives a root branch length, which is not accounted for by the above generative processes. We use the mean of the PC64 and VRC01 root branch lengths ( $\approx 0.01759$ ) as the default root branch length for simulation and also assign simulated trees root branch lengths of 0.1 to validate inference in settings with long periods of shared mutational history. For each simulated tree, we draw a naive sequence at the root according to the **partis** prior distribution that, as mentioned previously, models the V(D)J rearrangement process. We use the default settings in **partis** for human heavy chain data, since this represents the regime of most common interest.

To simulate DNA sequence data on the selected trees given the **partis**-generated naive sequences, we use the simulator in the **samm** package (Feng et al., 2019). A complex collection of enzymes introduces mutations to affinity-maturing sequences in a random pattern that is known to be highly sensitive to the sequence motif (i.e. the subsequence of DNA bases surrounding the mutating position) (Rogozin and Kolchanov, 1992; Dunn-Walters et al., 1998; Chahwan et al., 2012; Methot and Di Noia, 2017). The **samm** program estimates these motif mutabilities (i.e. the probability a position will mutate given the motif at that position) and substitution probabilities (i.e. the probability a position will mutate to a new base given the motif at that position) using a penalized Cox proportional hazards model and simulates sequence mutations given these inferred parameters. In our work, we use the default **samm** parameters inferred for human heavy chain sequences (Feng et al., 2019).

In our simulation experiments, we set the beta-splitting “balance” parameter  $\beta$  to  $-1.5$ ,  $-1.25$ , and  $-1$ ; the CF sequence count  $N_{\text{CF}}$  to 40 and 80; and the root branch length  $t_0$  to 0.01759 and 0.1. There are thus 12 different parameter combinations, and for each of these, we simulate 15 trees for a grand total of 180 simulated trees. While 180 simulated trees does not seem like a large number of Monte Carlo replicates, we are forced to constrain the number of trees on which we perform inference because the **ARPP** program has to be run by hand via a GUI. Using these parameter settings allows for the simulation of both balanced and imbalanced trees.



## Naive Sequence Validations

For each of the 180 simulated trees, we validate the accuracy of the naive sequence estimates generated from the **linearham**, **partis**, and **ARPP** programs by computing the hamming distances (i.e. the number of mismatched characters between two equal-length strings) between the estimates and the true naive sequence. We did this for the inferred DNA sequences directly and also for those same sequences translated to amino acids after inference was complete. The **partis** and **ARPP** packages provide naive sequence point estimates, but **linearham** samples naive sequences from a posterior distribution so we take the naive sequence with the highest posterior probability as the **linearham** “point estimate”.

We summarize the hamming distance results for all 180 simulated trees described above and plot this performance metric against the corresponding values of tree imbalance (Figure 3); for reference, we plot the tree imbalance values for the PC64 and VRC01 trees as well. The performance of **partis** clearly worsens as trees become more imbalanced, which makes sense given that **partis** assumes a star-tree configuration for clonal family evolution. The **linearham** and **ARPP** programs provide accurate naive sequence estimates and perform similarly across the observed tree imbalance spectrum, which is not too surprising because they both incorporate phylogenetic information into their estimates. We also show

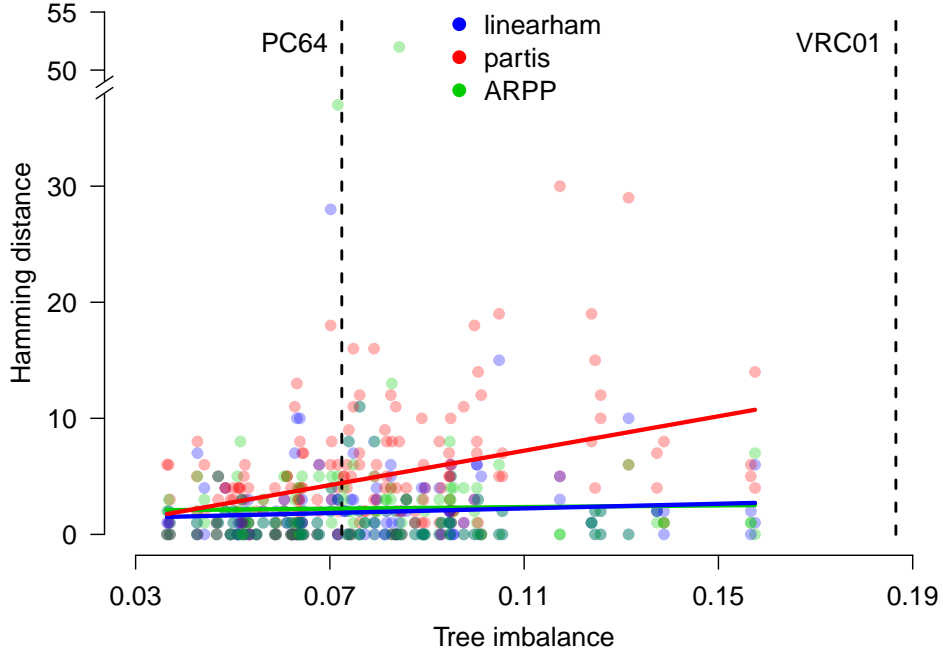


Figure 3: The hamming distances between the simulated naive DNA sequences and their corresponding **linearham**, **partis**, and **ARPP** estimates versus the tree imbalance values of the simulated trees. Linear regression lines are superimposed for each method to indicate how the results vary as trees get more imbalanced. For reference, we plot the tree imbalance values for the PC64 and VRC01 trees.

a summary of these results, split into values for the full-sequence and CDR3 regions, as well as the DNA/amino-acid sequence types (Table 1). Table 1 suggests that **linearham** and **ARPP** both perform better than **partis** does by approximately 2-3 nucleotides (1.6-1.8 amino acids) in the whole sequence and in the CDR3 region. Furthermore, it seems **linearham** and **ARPP** perform similarly across all the different settings in Table 1.

We now present the mean hamming distance results, averaging over all trees generated under the different simulation parameter settings. Specifically, we average over the trees that were simulated using beta-splitting “balance” parameter values  $\beta = -1, -1.25, -1.5$ , CF sequence counts  $N_{CF} = 40, 80$ , and root branch lengths  $t_0 = 0.01759, 0.1$ . The mean hamming distance values seem to increase slightly for **linearham** and **ARPP** and considerably for **partis** as we go from  $\beta = -1, -1.25$  to  $\beta = -1.5$  (Table S1), which is also suggested in Figure 3. The performance of **linearham** and **partis** deteriorates as  $N_{CF}$  goes from 40 to 80, while the opposite result is true for **ARPP** (Table S2). Despite this, **linearham** is still better than **ARPP** at predicting the whole naive sequence. As the root branch length  $t_0$  increases from 0.01759 to 0.1, the mean hamming distances increase substantially for all three programs (Table S3), which is

Sequence Region	Program	Sequence Type	
		DNA	Amino-Acid
Full-sequence	<b>linearham</b>	1.92	1.17
		(3.14)	(1.76)
	<b>partis</b>	4.81	3.02
		(4.91)	(2.84)
	<b>ARPP</b>	2.24	1.39
		(5.06)	(3.32)
CDR3-only	<b>linearham</b>	1.66	1.08
		(2.89)	(1.69)
	<b>partis</b>	4.17	2.69
		(3.70)	(2.29)
	<b>ARPP</b>	1.27	0.856
		(1.78)	(1.18)

Table 1: Mean hamming distances between the simulated naive sequences and their corresponding estimates, where the hamming distances are averaged over all 180 simulated trees. Results are provided for the **linearham**, **partis**, and **ARPP** programs; the full-sequence and CDR3 regions; and the DNA/amino-acid sequence types. Standard errors are also presented in parentheses.

intuitive because we are essentially introducing a higher number of mutations that are shared across all clonal sequences. While the **linearham** naive sequence validations did not take advantage of the full naive sequence posterior distribution, we demonstrate the usefulness of accounting for phylogenetic uncertainty in our ancestral sequence validations described below.

## Intermediate Ancestral Sequence Validations

Our ancestral sequence validation experiments are centered around accurate inference of particular root-to-tip ancestral sequence lineages of interest because immunologists frequently use ASR to estimate the mutational pathways associated with antibody development (Doria-Rose et al., 2014; Simonich et al., 2019). For each of our simulated trees, we determine the root-to-tip ancestral lineage of interest by identifying the tip sequence that is farthest from the naive sequence in terms of branch length distance.

We quantify the results of our ancestral sequence validation by treating it as a machine learning classification problem: do the posterior probabilities aid us in deciding if the ancestral lineage sequences are correct? In all our experiments, we measure classification performance by recording the positive predictive value (i.e. the fraction of sequences in the *ancestral lineage prediction set* that are on the *true ancestral lineage*) and the true positive rate (i.e. the fraction of sequences on the *true ancestral lineage* that are in the *ancestral lineage prediction set*). The “predicted classification” of these sequences is obtained by applying a decision boundary  $\rho \in \{0.25, 0.5, 0.75\}$  to the posterior probability. Thus, for example, if  $\rho = 0.75$  and a given ancestral sequence is on the true ancestral lineage and has posterior probability 0.8, it is considered to be a “true positive” prediction. This analysis is straightforward for the **linearham** and **RevBayes** programs, that do estimate posterior probabilities. We define **dnaml** “posterior probabilities” to be either 0 or 1 depending on whether a lineage sequence is either outside or inside of the **dnaml**-inferred set of most probable reconstructed lineage sequences.

We report the positive predictive values and the true positive rates for all 180 simulated trees and plot these performance metrics against the corresponding values of tree imbalance (Figure 4); for reference, we plot the tree imbalance values for the PC64 and VRC01 trees as well. Notice that the superimposed linear regression lines have negative slopes close to 0, which suggests that ancestral lineage inference is not clearly sensitive to the “balance” of the tree. We also display the mean positive predictive values and mean true positive rates aggregated over all 180 trees for the same three programs, the different decision boundaries  $\rho \in \{0.25, 0.5, 0.75\}$ , and the DNA/amino-acid sequence types (Table 2). Table 2 indicates that **linearham** performs better than **RevBayes** in every setting, indicating that accounting for naive rearrangement uncertainty in our posterior distribution rather than conditioning on the **partis**-inferred naive sequence leads to more accurate ancestral lineage sequence estimates. At the lowest decision boundary  $\rho = 0.25$ , **linearham** obtains slightly better positive predictive values and true positive rates than either **RevBayes** or **dnaml** does (Table 2). As the decision boundary  $\rho$  increases, **linearham** and **RevBayes** achieve higher positive predictive values at the expense of lower true positive rates, which makes sense because high values of  $\rho$  imply that only lineage sequences with high posterior probabilities



are predicted to be on the true lineage. In addition, note that the increases in positive predictive values are greater than the decreases in true positive rates for **linearham** as  $\rho$  increases. Of course, **dnaml** obtains the same positive predictive values and true positive rates regardless of  $\rho$  because its “posterior probabilities” are either 0 or 1.

We also present the mean positive predictive values and mean true positive rates for decision boundary  $\rho = 0.5$ , averaging over all trees generated under the different simulation parameter settings. The validation performance of all the programs seems to decline, albeit slightly, from  $\beta = -1$  to  $\beta = -1.5$  (Table S4), which makes sense given the trends in Figure 4. For the most part, mean positive predictive values and mean true positive rates increase for **linearham**, **RevBayes**, and **dnaml** as  $N_{CF}$  goes from 40 to 80 (Table S5), which seems surprising but for larger CFs, the root-to-tip lineage becomes larger and may explain this observed pattern. As the root branch length  $t_0$  increases from 0.01759 to 0.1, the ancestral lineage validation results worsen considerably for all three programs (Table S6), which is logical because, as we stated above, we are essentially introducing a higher number of mutations shared across all clonal sequences when we utilize longer root branch lengths in our simulations.

These results help demonstrate why Bayesian ancestral lineage inference should be favored over likelihood-based approaches to intermediate lineage inference as the Bayesian posterior probabilities quantify the uncertainty in our sequence estimates. In a real-life experimental setting, the decision boundary  $\rho$  should be chosen based on the desired level of positive predictive values or true positive rates and our analysis provides some insight into the mapping between  $\rho$  and these classification metrics. In practice, immunologists are probably more interested in controlling positive predictive values than true positive rates because synthesizing computationally-inferred lineage sequences and testing their binding and neutralization abilities is a laborious and expensive endeavor. Thus, knowing the approximate fraction of inferred intermediate lineage sequences that are on the true lineage is of the utmost importance to an immunologist.

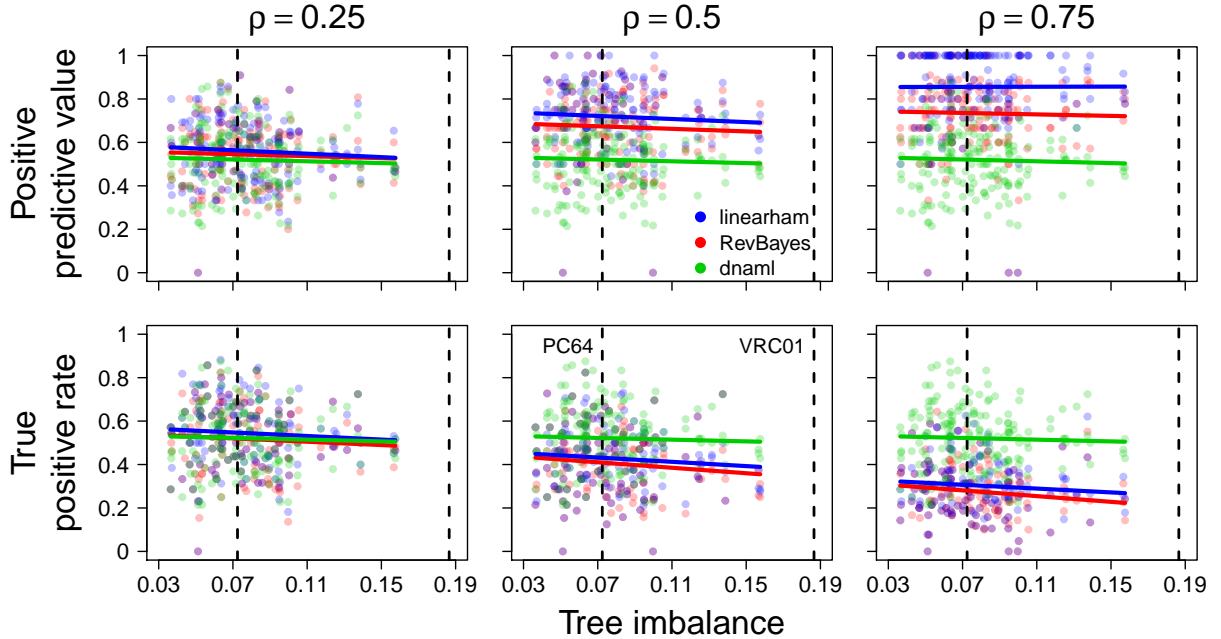


Figure 4: The positive predictive values and the true positive rates versus the tree imbalance values of the simulated trees, stratified by decision boundary  $\rho$ . Positive predictive values and true positive rates are computed on the DNA sequences and for the **linearham**, **RevBayes**, and **dnaml** programs. Linear regression lines are superimposed for each package to indicate how the results vary as trees get more imbalanced. For reference, we plot the tree imbalance values for the PC64 and VRC01 trees (vertical dashed lines).

Performance Metric	Program	Sequence Type					
		DNA			Amino-Acid		
		$\rho = 0.25$	$\rho = 0.5$	$\rho = 0.75$	$\rho = 0.25$	$\rho = 0.5$	$\rho = 0.75$
Positive predictive value	<b>linearham</b>	0.561	0.719	0.856	0.613	0.758	0.858
		(0.139)	(0.157)	(0.176)	(0.131)	(0.145)	(0.146)
	<b>RevBayes</b>	0.544	0.672	0.734	0.590	0.713	0.774
		(0.141)	(0.160)	(0.166)	(0.134)	(0.145)	(0.139)
	<b>dnaml</b>	0.520	0.520	0.520	0.590	0.590	0.590
		(0.139)	(0.139)	(0.139)	(0.165)	(0.165)	(0.165)
True positive rate	<b>linearham</b>	0.545	0.428	0.304	0.640	0.533	0.398
		(0.146)	(0.147)	(0.125)	(0.139)	(0.144)	(0.138)
	<b>RevBayes</b>	0.517	0.406	0.276	0.606	0.505	0.370
		(0.147)	(0.144)	(0.116)	(0.140)	(0.144)	(0.134)
	<b>dnaml</b>	0.521	0.521	0.521	0.584	0.584	0.584
		(0.142)	(0.142)	(0.142)	(0.166)	(0.166)	(0.166)

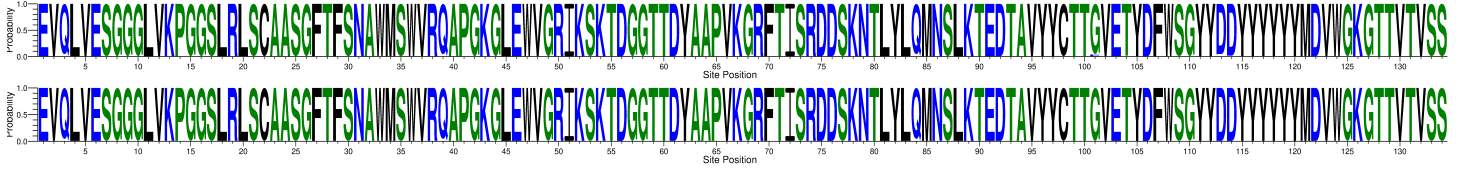
Table 2: Mean positive predictive values and mean true positive rates, averaged over all 180 simulated trees. Results are provided for the **linearham**, **RevBayes**, and **dnaml** programs; the different decision boundaries  $\rho \in \{0.25, 0.5, 0.75\}$ ; and the DNA/amino-acid sequence types. Standard errors are also presented in parentheses.

## PC64/VRC01 Ancestral Lineage Analysis

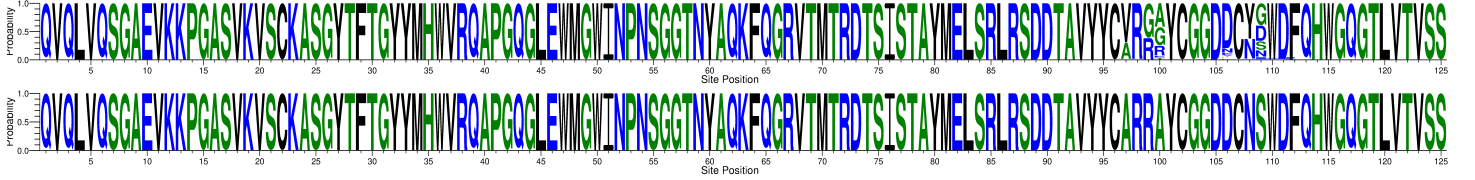
We illustrate the capabilities of **linearham** and **ARPP/dnaml** on real-world datasets by applying the three methods to subsets of the PC64 and VRC01 datasets. The PC64 dataset contains a set of clonal sequences with multiple HIV-binding bNAbs that results from a longitudinal study over 46 months on an African donor (i.e. donor PC64) within the International AIDS Vaccine Initiative Protocol C cohort (Landais et al., 2016). Our VRC01 clonal family dataset also originates from an HIV-infected donor and contains many bNAb sequences that are part of a well-known class of HIV-binding antibodies (i.e. the VRC01 class) (Wu et al., 2011; West et al., 2012; Zhou et al., 2013; Wu et al., 2015). The tip sequences of interest for the PC64 and VRC01 datasets are chosen to be PCT64-35M and NIH45-46, respectively, which are both monoclonal antibody sequences that have accumulated a large amount of SHM. We use 100 sequences from the PC64 CF dataset using a pruning strategy discussed in (Simonich et al., 2019) and trim the VRC01 CF dataset to 268 sequences using the **cd-hit** sequence clustering program (Li and Godzik, 2006) with a 95% sequence identity cutoff. We perform inference on these subsetted datasets using the same settings as described for our simulation experiments.

The PC64 amino acid naive sequence posterior probability logos suggest that there is little uncertainty in the naive sequence reconstruction (Figure 5a) and, in fact, the most probable **linearham** amino acid naive sequence has a probability of approximately 0.92. However, the VRC01 naive sequence seems to have considerable uncertainty in the CDR3 region (Figure 5b), which shows that properly modeling ancestral sequence and phylogenetic uncertainty is important for real-world datasets with highly-mutated sequences. The most probable **linearham** VRC01 naive sequence estimate has a posterior probability approximately equal to 0.036. It is important to note that the VRC01 CF sequences were first collected 5 years after the diagnosis of the associated HIV-1 infection, whereas the PC64 CF sequences contain samples drawn from the corresponding donor as early as 4 months post infection. This indicates that the VRC01 naive sequence reconstruction inherently has more uncertainty because there are not any early time-point samples in the corresponding dataset. Furthermore, the most probable **linearham** VRC01 naive sequence amino acids do not perfectly match the corresponding **ARPP**-inferred residues (Figure 5b). In total, these results suggest that in the absence of uncertainty, **linearham** and **ARPP** produce similar naive sequence reconstructions. However, when there is a significant amount of naive sequence uncertainty, **linearham**, unlike **ARPP**, provides alternate hypotheses that should be considered along with corresponding uncertainty estimates.

Our **linearham** analysis demonstrates that there are many probable naive-to-tip sequence paths (Figure 6), suggesting that intermediate ancestral sequences have high levels of uncertainty; we use 0.04 probability cutoffs in each **linearham**-inferred lineage graphic. In particular, the **linearham** lineage diagram for the PC64 dataset (Figure 6a) shows many possible routes of evolution from the different naive sequences to the PCT64-35M mature sequence and displays confidence values via posterior probabilities. The VRC01 lineage graphic in Figure 6b does not display connections between any naive sequences



(a) PC64 naive sequence posterior probability logos



(b) VRC01 naive sequence posterior probability logos

Figure 5: The **linearham**-inferred (top) and **ARPP**-inferred (bottom) amino acid naive sequence posterior probability logos for the pruned PC64 dataset of 100 sequences and the trimmed VRC01 alignment of 268 sequences.

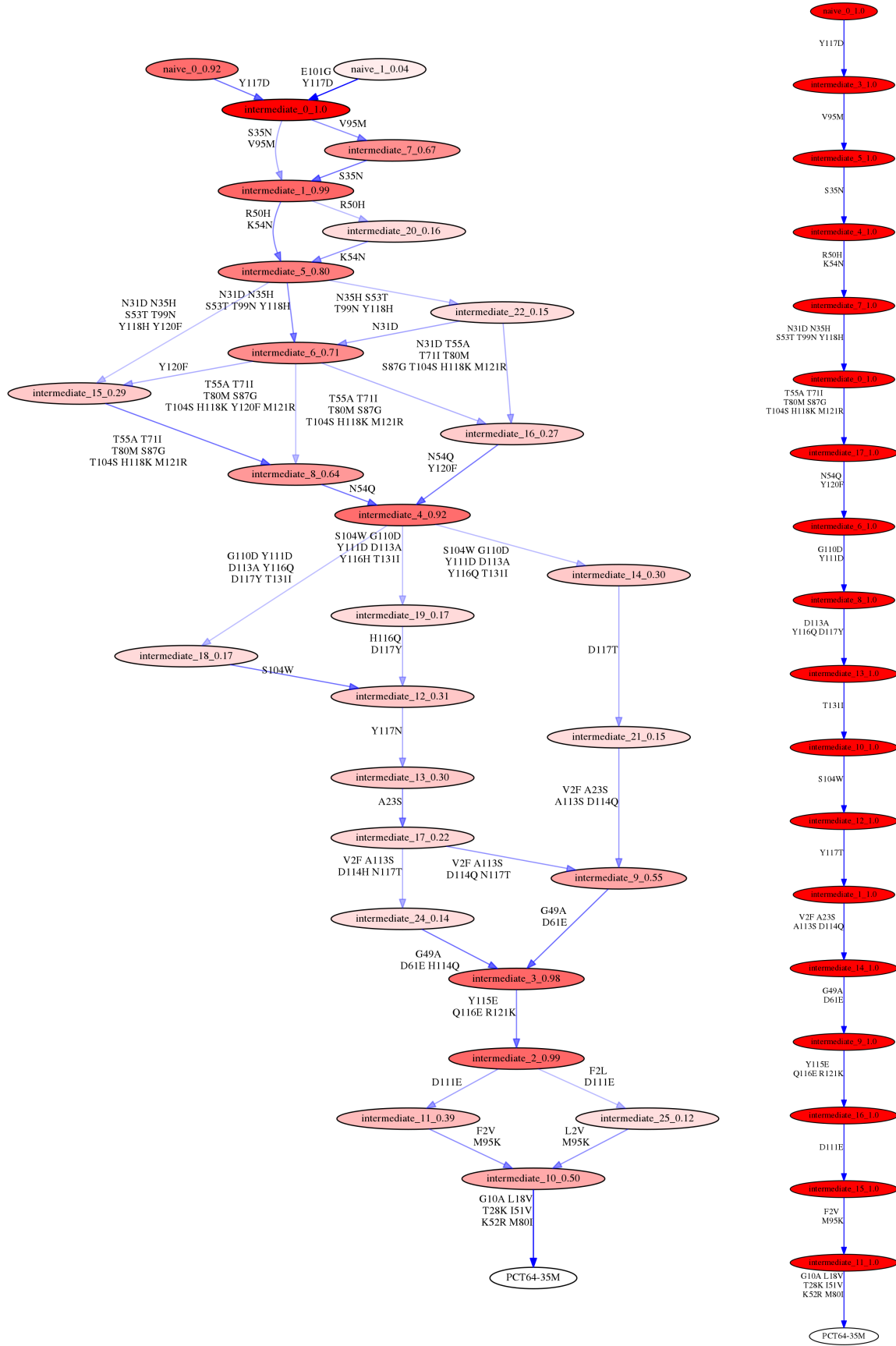
and intermediate ancestral sequences, which reflects the fact that naive sequence inference is extremely challenging on this dataset. In summary, there is considerable uncertainty in naive-to-tip mutational trajectory inference in real-world BCR datasets. This finding contradicts the assertions of [Kepler \(2013\)](#), who states that ancestral sequence and phylogenetic uncertainty is unimportant, and proceeds with a single point estimate.

## Discussion

In this paper, we introduce a novel Bayesian approach to CF phylogenetic inference that is based on a phylo-HMM. Our phylo-HMM posterior sampling methodology not only allows for easy quantification of phylogenetic and ancestral sequence uncertainty but also models the V(D)J recombination process as an informative prior on the root sequence. Specifically, our phylo-HMM models both the naive rearrangement and SHM processes by using a hidden state discrete-time Markov model on naive sequences that explicitly incorporates V(D)J rearrangement information and an emission distribution generating the clonal sequences conditional on the naive sequence that is based on phylogenetic likelihoods. We show that our inference procedure, implemented in the software package **linearham**, provides higher-quality naive sequence and ancestral sequence estimates compared to those obtained under current state-of-the-art methods and augments these estimates with relative confidence values by reporting the associated posterior probabilities.

From our simulation experiments, we see that the **partis** naive sequence estimates get substantially worse as trees get more imbalanced. This is in contrast to those of **linearham** and **ARPP**, which is intuitive because **partis** assumes a star-tree configuration whereas **linearham** and **ARPP** leverage phylogenetic models. The **linearham** and **ARPP** programs perform similarly in our naive sequence validations regarding their most highly supported inferences, but only **linearham** is capable of characterizing naive sequence uncertainty. This is important because, as we see in our VRC01 analysis (and in contrast to the findings of [Kepler \(2013\)](#)), real datasets of significant practical importance can have very large uncertainties associated with their naive sequence estimates. In addition, we demonstrate that **linearham** ancestral lineage inference performs better, via mean positive predictive values and mean true positive rates, than either **RevBayes** or **dnaml** does at the lowest decision boundary  $\rho = 0.25$ , which suggests that accounting for naive sequence, phylogenetic, and ancestral sequence uncertainty does lead to slightly improved ASR performance. Furthermore, our Bayesian ASR results indicate that  $\rho$  can be chosen according to a prespecified trade-off between positive predictive values and true positive rates, which is an important consideration for immunologists looking to synthesize computationally-inferred lineage sequences.

Based on all the evidence in this manuscript, we recommend using **linearham** to infer and visualize naive-to-mature mutational pathways in CFs harboring antibodies of interest. Bayesian phylogenetic inference has already been shown to be useful for identifying different possible routes of evolution from a fixed naive sequence to a bNAb with relative confidence values ([Simonich et al., 2019](#)) and we believe



(a) PC64 posterior probability lineage inference



(b) VRC01 posterior probability lineage inference

Figure 6: The `linearham`-inferred (left) and `dnaml`-inferred (right) naive-to-tip amino acid sequence trajectories for the pruned PC64 dataset of 100 sequences and the trimmed VRC01 alignment of 268 sequences. The tip sequences of interest for the PC64 and VRC01 datasets are chosen to be PCT64-35M and NIH45-46, respectively, and we use 0.04 probability cutoffs for these lineage graphics. The nodes correspond to unique ancestral sequences filled with red color, where the opacity is proportional to the posterior probability of the associated sequence. Each node has a label that denotes whether the associated sequence is a naive or intermediate ancestral sequence, the posterior probability rank of the sequence among all sampled naive or intermediate ancestral sequences, and the sequence-specific posterior probability itself. The directed edges connecting nodes represent ancestral sequence transitions, are shaded blue with an opacity proportional to the posterior probability of the associated sequence transition, and are annotated with the site-specific mutations between the two sequences.

our Bayesian phylo-HMM analysis pipeline in **linearham** can be used to not only infer similarly-styled maturation pathways but also visualize the uncertainties inherent in the naive sequence and ancestral sequence estimates. From a practical standpoint, these different possible ancestral lineages allow immunologists to generate many different intermediate antibody candidates that bind and/or neutralize HIV.

Our Bayesian phylo-HMM inference procedure admits a number of possible future extensions to enhance the effectiveness of our technique. One drawback of our method is that it does not sample the parameters of  $p(\mathbf{Y}_{\text{naive}})$  and uses **partis** (and its star-tree assumption) to estimate them. Ideally, our Bayesian inference procedure would jointly sample all the model parameters, but currently this is not practically feasible. In addition, our method implicitly assumes the CFs have been obtained from **partis**, which uses the star-tree assumption to cluster repertoire sequences. It may be possible to, as [Ralph and Matsen IV \(2016b\)](#) state, incorporate the phylo-HMM in the CF clustering procedure within **partis** to obtain higher-quality CFs, but in our current Bayesian implementation that would be computationally costly.

## Acknowledgments

We thank Arman Bilge for many stimulating conversations about our phylo-HMM emission likelihood calculation, Andy Magee for answering all our questions about **RevBayes**, and Jean Feng for adapting the sequence simulator in **samm** to be used in our simulations. We also want to thank Chaim Schramm for providing us with the VRC01 sequence dataset used in ([Wu et al., 2015](#)), Bryan Briney for sending us the PC64 multiple sequence alignment used in ([Landais et al., 2017](#)), Tyler Starr for creating the trimmed 268-sequence VRC01 dataset described above, and Laura Doepker for providing us with other experimental datasets to test **linearham** on. This research was supported by NIH grants R01-GM113246, R01-AI120961, R01-AI138709, and U19-AI117891 as well as National Science Foundation grants CISE-1561334 and CISE-1564137. The research of Frederick Matsen was supported in part by a Faculty Scholar grant from the Howard Hughes Medical Institute and the Simons Foundation. Amrit Dhar was supported by an NSF IGERT DGE-1258485 fellowship.

## References

- Aldous, D. (1996). Probability distributions on cladograms. In *Random Discrete Structures*, pages 1–18. Springer.
- Andrieu, C., Doucet, A., and Holenstein, R. (2010). Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 72(3):269–342.
- Chahwan, R., Edelmann, W., Scharff, M. D., and Roa, S. (2012). AIDing antibody diversity by error-prone mismatch repair. In *Seminars in Immunology*, volume 24, pages 293–300. Elsevier.
- Crooks, G. E., Hon, G., Chandonia, J.-M., and Brenner, S. E. (2004). WebLogo: a sequence logo generator. *Genome Research*, 14(6):1188–1190.
- Doria-Rose, N. A., Bhiman, J. N., Roark, R. S., Schramm, C. A., Gorman, J., Chuang, G.-Y., Pancera, M., Cale, E. M., Ernandes, M. J., Louder, M. K., et al. (2016). New member of the V1V2-directed CAP256-VRC26 lineage that shows increased breadth and exceptional potency. *Journal of Virology*, 90(1):76.
- Doria-Rose, N. A., Schramm, C. A., Gorman, J., Moore, P. L., Bhiman, J. N., DeKosky, B. J., Ernandes, M. J., Georgiev, I. S., Kim, H. J., Pancera, M., et al. (2014). Developmental pathway for potent V1V2-directed HIV-neutralizing antibodies. *Nature*, 509(7498):55.
- Dunn-Walters, D. K., Dogan, A., Boursier, L., MacDonald, C. M., and Spencer, J. (1998). Base-specific sequences that bias somatic hypermutation deduced by analysis of out-of-frame human IgVH genes. *The Journal of Immunology*, 160(5):2360–2364.
- Elhanati, Y., Sethna, Z., Marcou, Q., Callan Jr, C. G., Mora, T., and Walczak, A. M. (2015). Inferring processes underlying B-cell repertoire diversity. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 370(1676).



- Felsenstein, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. Journal of Molecular Evolution, 17(6):368–376.
- Felsenstein, J. (2004). Inferring Phylogenies. Sinauer Associates, Sunderland, Massachusetts.
- Felsenstein, J. (2005). PHYLIP (Phylogeny Inference Package) version 3.6. Distributed by the author. Department of Genome Sciences, University of Washington, Seattle.
- Feng, J., Shaw, D. A., Minin, V. N., Simon, N., and Matsen IV, F. A. (2019). Survival analysis of DNA mutation motifs with penalized proportional hazards. Annals of Applied Statistics, 13(2):1268–1294.
- Flouri, T., Izquierdo-Carrasco, F., Darriba, D., Aberer, A., Nguyen, L., Minh, B., Von Haeseler, A., and Stamatakis, A. (2014). The phylogenetic likelihood library. Systematic Biology, 64(2):356–362.
- Gansner, E. R. and North, S. C. (2000). An open graph visualization system and its applications to software engineering. Software: Practice and Experience, 30(11):1203–1233.
- Gelman, A. and Meng, X.-L. (2004). Applied Bayesian Modeling and Causal Inference from Incomplete-Data Perspectives. John Wiley & Sons.
- Gelman, A., Stern, H. S., Carlin, J. B., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2013). Bayesian Data Analysis. Chapman and Hall/CRC.
- Gong, L. I., Suchard, M. A., and Bloom, J. D. (2013). Stability-mediated epistasis constrains the evolution of an influenza protein. eLife, 2:e00631.
- Gordon, N. J., Salmond, D. J., and Smith, A. F. (1993). Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In IEE Proceedings F (Radar and Signal Processing), volume 140, pages 107–113. IET.
- Hanson-Smith, V., Kolaczowski, B., and Thornton, J. W. (2010). Robustness of ancestral sequence reconstruction to phylogenetic uncertainty. Molecular Biology and Evolution, 27(9):1988–1999.
- Hoehn, K. B., Lunter, G., and Pybus, O. G. (2017). A phylogenetic codon substitution model for antibody lineages. Genetics, 206(1):417–427.
- Hoehn, K. B., Vander Heiden, J. A., Zhou, J. Q., Lunter, G., Pybus, O. G., and Kleinstein, S. (2019). Repertoire-wide phylogenetic models of B cell molecular evolution reveal evolutionary signatures of aging and vaccination. bioRxiv, page 558825.
- Höhna, S., Landis, M. J., Heath, T. A., Boussau, B., Lartillot, N., Moore, B. R., Huelsenbeck, J. P., and Ronquist, F. (2016). RevBayes: Bayesian phylogenetic inference using graphical models and an interactive model-specification language. Systematic Biology, 65(4):726–736.
- Kepler, T. B. (2013). Reconstructing a B-cell clonal lineage. I. statistical inference of unobserved ancestors. F1000Research, 2(103).
- Landais, E., Huang, X., Havenar-Daughton, C., Murrell, B., Price, M. A., Wickramasinghe, L., Ramos, A., Bian, C. B., Simek, M., Allen, S., et al. (2016). Broadly neutralizing antibody responses in a large longitudinal sub-saharan HIV primary infection cohort. PLoS Pathogens, 12(1):e1005369.
- Landais, E., Murrell, B., Briney, B., Murrell, S., Rantalainen, K., Berndsen, Z. T., Ramos, A., Wickramasinghe, L., Smith, M. L., Eren, K., et al. (2017). HIV envelope glycoform heterogeneity and localized diversity govern the initiation and maturation of a V2 apex broadly neutralizing antibody lineage. Immunity, 47(5):990–1003.
- Lauritzen, S. L. (1996). Graphical Models, volume 17. Clarendon Press.
- Li, W. and Godzik, A. (2006). Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. Bioinformatics, 22(13):1658–1659.
- Liao, H.-X., Lynch, R., Zhou, T., Gao, F., Alam, S. M., Boyd, S. D., Fire, A. Z., Roskin, K. M., Schramm, C. A., Zhang, Z., et al. (2013). Co-evolution of a broadly neutralizing HIV-1 antibody and founder virus. Nature, 496(7446):469.



- Mascola, J. R. and Haynes, B. F. (2013). HIV-1 neutralizing antibodies: understanding nature’s pathways. Immunological Reviews, 254(1):225–244.
- Method, S. and Di Noia, J. (2017). Molecular mechanisms of somatic hypermutation and class switch recombination. In Advances in Immunology, volume 133, pages 37–87. Elsevier.
- Nielsen, R. (2002). Mapping mutations on phylogenies. Systematic Biology, 51(5):729–739.
- Price, M. N., Dehal, P. S., and Arkin, A. P. (2009). FastTree: computing large minimum evolution trees with profiles instead of a distance matrix. Molecular Biology and Evolution, 26(7):1641–1650.
- Price, M. N., Dehal, P. S., and Arkin, A. P. (2010). FastTree 2—approximately maximum-likelihood trees for large alignments. PLoS One, 5(3):e9490.
- Rabiner, L. R. (1986). An introduction to hidden Markov models. IEEE ASSP Magazine, 3(1):4–16.
- Ralph, D. K. and Matsen IV, F. A. (2016a). Consistency of VDJ rearrangement and substitution parameters enables accurate B cell receptor sequence annotation. PLoS Computational Biology, 12(1):e1004409.
- Ralph, D. K. and Matsen IV, F. A. (2016b). Likelihood-based inference of B cell clonal families. PLoS Computational Biology, 12(10):e1005086.
- Rogozin, I. B. and Kolchanov, N. A. (1992). Somatic hypermutagenesis in immunoglobulin genes: II. influence of neighbouring base sequences on mutagenesis. Biochimica et Biophysica Acta (BBA)-Gene Structure and Expression, 1171(1):11–18.
- Rubin, D. B. (1987). The calculation of posterior distributions by data augmentation: Comment: A non-iterative sampling/importance resampling alternative to the data augmentation algorithm for creating a few imputations when fractions of missing information are modest: The SIR algorithm. Journal of the American Statistical Association, 82(398):543–546.
- Scott, S. L. (2002). Bayesian methods for hidden Markov models: Recursive computing in the 21st century. Journal of the American Statistical Association, 97(457):337–351.
- Siepel, A. and Haussler, D. (2005). Phylogenetic hidden Markov models. In Statistical Methods in Molecular Evolution, pages 325–351. Springer.
- Simonich, C. A., Doepker, L., Ralph, D., Williams, J. A., Dhar, A., Yaffe, Z., Gentles, L., Small, C. T., Oliver, B., Vigdorovich, V., et al. (2019). Kappa chain maturation helps drive rapid development of an infant HIV-1 broadly neutralizing antibody lineage. Nature Communications, 10(1):2190.
- Skare, Ø., Bølviken, E., and Holden, L. (2003). Improved sampling-importance resampling and reduced bias importance sampling. Scandinavian Journal of Statistics, 30(4):719–737.
- Smith, A. F. and Gelfand, A. E. (1992). Bayesian statistics without tears: a sampling–resampling perspective. The American Statistician, 46(2):84–88.
- Stamatatos, L., Pancera, M., and McGuire, A. T. (2017). Germline-targeting immunogens. Immunological Reviews, 275(1):203–216.
- Tavaré, S. (1986). Some probabilistic and statistical problems in the analysis of DNA sequences. Lectures on Mathematics in the Life Sciences, 17:57–86.
- Watson, C. T., Glanville, J., and Marasco, W. A. (2017). The individual and population genetics of antibody immunity. Trends Immunol., 38(7):459–470.
- West, A. P., Diskin, R., Nussenzweig, M. C., and Bjorkman, P. J. (2012). Structural basis for germ-line gene usage of a potent class of antibodies targeting the CD4-binding site of HIV-1 gp120. Proceedings of the National Academy of Sciences, 109(30):E2083–E2090.
- Wu, X., Zhang, Z., Schramm, C. A., Joyce, M. G., Do Kwon, Y., Zhou, T., Sheng, Z., Zhang, B., O’Dell, S., McKee, K., et al. (2015). Maturation and diversity of the VRC01-antibody lineage over 15 years of chronic HIV-1 infection. Cell, 161(3):470–485.

- Wu, X., Zhou, T., Zhu, J., Zhang, B., Georgiev, I., Wang, C., Chen, X., Longo, N. S., Louder, M., McKee, K., et al. (2011). Focused evolution of HIV-1 neutralizing antibodies revealed by structures and deep sequencing. Science, 333(6049):1593–1602.
- Yaari, G., Vander Heiden, J., Uduman, M., Gadala-Maria, D., Gupta, N., Stern, J., O'Connor, K., Hafler, D., Laserson, U., Vigneault, F., and Kleinstein, S. (2013). Models of somatic hypermutation targeting and substitution based on synonymous mutations from high-throughput immunoglobulin sequencing data. Frontiers in Immunology, 4:358.
- Yang, Z. (1994). Maximum likelihood phylogenetic estimation from DNA sequences with variable rates over sites: Approximate methods. Journal of Molecular Evolution, 39(3):306–314.
- Yang, Z. (1996). Among-site rate variation and its impact on phylogenetic analyses. Trends in Ecology & Evolution, 11(9):367–372.
- Zhou, T., Zhu, J., Wu, X., Moquin, S., Zhang, B., Acharya, P., Georgiev, I. S., Altae-Tran, H. R., Chuang, G.-Y., Joyce, M. G., et al. (2013). Multidonor analysis reveals structural elements, genetic determinants, and maturation pathway for HIV-1 neutralization by VRC01-class antibodies. Immunity, 39(2):245–258.

## Supplementary Figures and Tables

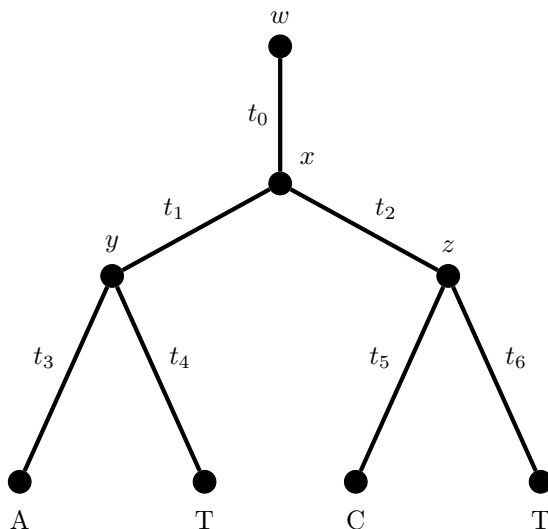


Figure S1: An example phylogenetic tree. Letters  $x, y, z, w$  represent the unobserved internal node states where  $w$  is associated with the root node,  $t_0$  defines the root branch length, and  $(t_0, t_1, t_2, \dots, t_6)$  denotes the entire vector of branch lengths. Given this tree topology and set of branch lengths, we can calculate the likelihood of observing the nucleotide vector  $(A, T, C, T)$  by marginalizing probabilities over the unobserved states  $x, y, z, w$ .

Sequence Region	Program	Sequence Type					
		DNA			Amino-Acid		
		$\beta = -1$	$\beta = -1.25$	$\beta = -1.5$	$\beta = -1$	$\beta = -1.25$	$\beta = -1.5$
Full-sequence	linearham	1.73	1.92	2.10	0.917	1.18	1.40
		(3.94)	(2.66)	(2.71)	(1.77)	(1.66)	(1.84)
	partis	3.85	3.82	6.77	2.48	2.42	4.15
		(3.26)	(3.85)	(6.48)	(1.90)	(2.41)	(3.61)
	ARPP	2.42	2.52	1.78	1.38	1.52	1.27
		(6.89)	(5.11)	(1.98)	(4.41)	(3.47)	(1.39)
CDR3-only	linearham	1.45	1.67	1.87	0.833	1.10	1.30
		(3.42)	(2.55)	(2.66)	(1.63)	(1.63)	(1.81)
	partis	3.53	3.47	5.50	2.38	2.25	3.45
		(2.84)	(3.45)	(4.34)	(1.81)	(2.20)	(2.64)
	ARPP	1.12	1.32	1.38	0.667	0.867	1.03
		(1.76)	(1.95)	(1.63)	(1.07)	(1.20)	(1.26)

Table S1: Mean hamming distances between the simulated naive sequences and their corresponding estimates, where the hamming distances are averaged over all trees generated under the different beta-splitting “balance” parameter value settings. Results are provided for the **linearham**, **partis**, and **ARPP** programs; the full-sequence and CDR3 regions; and the DNA/amino-acid sequence types. Standard errors are also presented in parentheses.

Sequence Region	Program	Sequence Type			
		DNA		Amino-Acid	
		$N_{CF} = 40$	$N_{CF} = 80$	$N_{CF} = 40$	$N_{CF} = 80$
Full-sequence	<b>linearham</b>	1.73	2.10	1.17	1.17
		(2.36)	(3.77)	(1.67)	(1.86)
	<b>partis</b>	4.12	5.50	2.72	3.31
		(3.58)	(5.90)	(2.40)	(3.20)
	<b>ARPP</b>	2.39	2.09	1.52	1.26
		(4.36)	(5.70)	(2.92)	(3.69)
CDR3-only	<b>linearham</b>	1.51	1.81	1.04	1.11
		(2.26)	(3.41)	(1.59)	(1.80)
	<b>partis</b>	3.84	4.49	2.57	2.82
		(3.38)	(3.98)	(2.29)	(2.30)
	<b>ARPP</b>	1.33	1.21	0.944	0.767
		(1.63)	(1.92)	(1.21)	(1.15)

Table S2: Table analogous to Table S1, but varying the CF sequence count  $N_{CF}$ .

Sequence Region	Program	Sequence Type			
		DNA		Amino-Acid	
		$t_0 = 0.01759$	$t_0 = 0.1$	$t_0 = 0.01759$	$t_0 = 0.1$
Full-sequence	<b>linearham</b>	0.744	3.09	0.456	1.88
		(1.27)	(3.94)	(0.889)	(2.10)
	<b>partis</b>	2.69	6.93	1.77	4.27
		(2.97)	(5.54)	(1.91)	(3.06)
	<b>ARPP</b>	0.80	3.68	0.467	2.31
		(1.26)	(6.76)	(0.737)	(4.46)
CDR3-only	<b>linearham</b>	0.644	2.68	0.433	1.72
		(1.17)	(3.65)	(0.875)	(2.04)
	<b>partis</b>	2.50	5.83	1.68	3.71
		(2.47)	(3.97)	(1.62)	(2.42)
	<b>ARPP</b>	0.444	2.10	0.311	1.40
		(0.751)	(2.10)	(0.574)	(1.37)

Table S3: Table analogous to Table S1, but varying the root branch length  $t_0$ .

Performance Metric	Program	Sequence Type					
		DNA			Amino-Acid		
		$\beta = -1$	$\beta = -1.25$	$\beta = -1.5$	$\beta = -1$	$\beta = -1.25$	$\beta = -1.5$
Positive predictive value	<b>linearham</b>	0.742	0.715	0.700	0.782	0.759	0.733
		(0.157)	(0.157)	(0.158)	(0.164)	(0.147)	(0.118)
	<b>RevBayes</b>	0.691	0.676	0.649	0.725	0.726	0.689
		(0.162)	(0.152)	(0.165)	(0.166)	(0.140)	(0.124)
	<b>dnaml</b>	0.519	0.528	0.512	0.575	0.605	0.590
		(0.141)	(0.149)	(0.129)	(0.179)	(0.173)	(0.141)
True positive rate	<b>linearham</b>	0.435	0.443	0.407	0.543	0.544	0.512
		(0.149)	(0.139)	(0.152)	(0.159)	(0.139)	(0.131)
	<b>RevBayes</b>	0.413	0.425	0.378	0.510	0.522	0.482
		(0.142)	(0.133)	(0.154)	(0.160)	(0.136)	(0.133)
	<b>dnaml</b>	0.521	0.531	0.512	0.566	0.600	0.586
		(0.144)	(0.152)	(0.132)	(0.182)	(0.166)	(0.148)

Table S4: Mean positive predictive values and mean true positive rates for decision boundary  $\rho = 0.5$ , where we average over all trees generated under the different beta-splitting “balance” parameter value settings. Results are provided for the **linearham**, **RevBayes**, and **dnaml** programs and the DNA/amino-acid sequence types. Standard errors are also presented in parentheses.

Performance Metric	Program	Sequence Type			
		DNA		Amino-Acid	
		$N_{CF} = 40$	$N_{CF} = 80$	$N_{CF} = 40$	$N_{CF} = 80$
Positive predictive value	<b>linearham</b>	0.714	0.725	0.759	0.756
		(0.165)	(0.150)	(0.166)	(0.121)
	<b>RevBayes</b>	0.654	0.690	0.709	0.717
		(0.171)	(0.147)	(0.168)	(0.117)
	<b>dnaml</b>	0.513	0.526	0.586	0.594
		(0.146)	(0.132)	(0.177)	(0.153)
True positive rate	<b>linearham</b>	0.430	0.427	0.527	0.539
		(0.153)	(0.141)	(0.159)	(0.127)
	<b>RevBayes</b>	0.403	0.408	0.501	0.509
		(0.151)	(0.137)	(0.158)	(0.129)
	<b>dnaml</b>	0.514	0.528	0.573	0.595
		(0.152)	(0.132)	(0.174)	(0.156)

Table S5: Table analogous to Table S4, but varying the CF sequence count  $N_{CF}$ .

Performance Metric	Program	Sequence Type			
		DNA		Amino-Acid	
		$t_0 = 0.01759$	$t_0 = 0.1$	$t_0 = 0.01759$	$t_0 = 0.1$
Positive predictive value	<b>linearham</b>	0.728	0.711	0.768	0.748
		(0.141)	(0.173)	(0.140)	(0.150)
	<b>RevBayes</b>	0.682	0.662	0.727	0.700
		(0.146)	(0.173)	(0.135)	(0.153)
	<b>dnaml</b>	0.537	0.503	0.616	0.564
		(0.133)	(0.144)	(0.154)	(0.172)
True positive rate	<b>linearham</b>	0.450	0.407	0.560	0.506
		(0.136)	(0.154)	(0.140)	(0.143)
	<b>RevBayes</b>	0.415	0.396	0.523	0.487
		(0.131)	(0.155)	(0.135)	(0.151)
	<b>dnaml</b>	0.536	0.506	0.609	0.559
		(0.136)	(0.147)	(0.157)	(0.171)

Table S6: Table analogous to Table S4, but varying the root branch length  $t_0$ .